

5

10

SYSTEM AND METHOD FOR DEBUGGING DISTRIBUTED SOFTWARE ENVIRONMENTS

Related Applications

15

This application is a continuation of U.S. Provisional Application No. 60/213,496 filed June 23, 2000, incorporated herein by reference.

Technical Field

20

The present invention relates to a system and method for debugging distributed software environments and, in particular to a coordination-centric approach in which a distributed software environment produces event traces to be analyzed by a debugging host.

Background of the Invention

25

A system design and programming methodology is most effective when it is closely integrated and coheres tightly with its corresponding debugging techniques. In distributed and embedded system methodologies, the relationship between debugging approaches and design methodologies has traditionally been one-sided in favor of the design and programming methodologies. Design and programming methodologies are typically developed without any consideration for the debugging techniques that will later be applied to software systems designed using that design and programming

30

methodology. While these typical debugging approaches attempt to exploit features provided by the design and programming methodologies, the debugging techniques will normally have little or no impact on what the design and programming features are in the first place. This lack of input from debugging approaches to design and programming methodologies serves to maintain the role of debugging as an afterthought, even though in a typical system design, debugging consumes a majority of the design time. The need remains for a design and programming methodology that reflects input from, and consideration of, potential debugging approaches in order to enhance the design and reduce the implementation time of software systems.

1. Packaging of Software Elements

Packaging refers to the set of interfaces a software element presents to other elements in a system. Software packaging has many forms in modern methodologies. Some examples are programming language procedure call interfaces (as with libraries), TCP/IP socket interfaces with scripting languages (as with mail and Web servers), and file formats. Several typical prior art packaging styles are described below, beginning with packaging techniques used in object-oriented programming languages and continuing with a description of more generalized approaches to packaging.

A. Object-Oriented Approaches to Packaging

One common packaging style is based on object-oriented programming languages and provides procedure-based (method-based) packaging for software elements (objects within this framework). These procedure-based packages allow polymorphism (in which several types of objects can have identical interfaces) through subtyping, and code sharing through inheritance (deriving a new class of objects from an already existing class of objects). In a typical object-oriented programming language, an object's interface is defined by the object's methods.

Object-oriented approaches are useful in designing concurrent systems (systems with task level parallelism and multiple processing resources?) because of the availability of active objects (objects with a thread of control). Some common,

concurrent object-oriented approaches are shown in actor languages and in concurrent Eiffel.

Early object-oriented approaches featured anonymity of objects through dynamic typechecking. This anonymity of objects meant that a first object did not
5 need to know anything about a second object in order to send a message to the second object. One unfortunate result of this anonymity of objects was that the second object could unexpectedly respond to the first object that the sent message was not understood, resulting in a lack of predictability, due to this disruption of system executions, for systems designed with this object-oriented approach.

10 Most modern object-oriented approaches opt to sacrifice the benefits flowing from anonymity of objects in order to facilitate stronger static typing (checking to ensure that objects will properly communicate with one another before actually executing the software system). The main result of stronger static typing is improved system predictability. However, an unfortunate result of sacrificing the anonymity of
15 objects is a tighter coupling between those objects, whereby each object must explicitly classify, and include knowledge about, other objects to which it sends messages. In modern object-oriented approaches the package (interface) has become indistinguishable from the object and the system in which the object is a part.

The need remains for a design and programming methodology that combines
20 the benefits of anonymity for the software elements with the benefits derived from strong static typing of system designs.

B. Other Approaches to Packaging

Other packaging approaches provide higher degrees of separation between software elements and their respective packages than does the packaging in
25 object-oriented systems. For example, the packages in event-based frameworks are interfaces with ports for transmitting and receiving events. These provide loose coupling for inter-element communication. However, in an event-based framework, a software designer must explicitly implement inter-element state coherence between software elements as communication between those software elements. This means

that a programmer must perform the error-prone task of designing, optimizing, implementing, and debugging a specialized communication protocol for each state coherence requirement in a particular software system.

The common object request broker architecture (CORBA) provides an interface description language (IDL) for building packages around software elements written in a variety of languages. These packages are remote procedure call (RPC) based and provide no support for coordinating state between elements. With flexible packaging, an element's package is implemented as a set of co-routines that can be adapted for use with applications through use of adapters with interfaces complementary to the interface for the software element. These adapters can be application-specific—used only when the elements are composed into a system.

The use of co-routines lets a designer specify transactions or sequences of events as part of an interface, rather than just as atomic events. Unfortunately, co-routines must be executed in lock-step, meaning a transition in one routine corresponds to a transition in the other co-routine. If there is an error in one or if an expected event is lost, the interface will fail because its context will be incorrect to recover from the lost event and the co-routines will be out of sync.

The need remains for a design and programming methodology that provides software packaging that supports the implementation of state coherence in distributed concurrent systems without packaging or interface failure when an error or an unexpected event occurs.

2. Approaches to Coordination

Coordination, within the context of this application, means the predetermined ways through which software components interact. In a broader sense, coordination refers to a methodology for composing concurrent components into a complete system. This use of the term coordination differs slightly from the use of the term in the parallelizing compiler literature, in which coordination refers to a technique for maintaining program-wide semantics for a sequential program decomposed into parallel subprograms.

Coordination languages are usually a class of tuple-space programming

The operators eval, inp, and readp are nonblocking versions of out, in, and read, respectively. They increase the expressive power of tuple-space languages.

15 Consider the case of eval, the nonblocking version of out. Instead of evaluating all arguments of the tuple before returning, it spawns a thread to evaluate them, creating, in effect, an active tuple (whereas tuples created by out are passive). As with out, when the computation is finished, the results are stored in a passive tuple and left in tuple space. Unlike out, however, the eval call returns immediately, so that several

20 active tuples can be left outstanding.

25

```

30  rpcCall(args) {
    out("RPCToServer", "Client", args...);
    in("Client", "ReturnFromServer", &returnValue);
    return returnValue;
}
Server:

```

```

...
while(true) {                                     /* C */
    in("RPCToServer", &returnAddress, args...);
    returnValue = functionCall(args);             /* C */
5      out(returnAddress, "ReturnFromServer", returnValue);
    }                                              /* C */

```

Although the implementation depicted in Listing 1 is a compact representation of an RPC protocol, the implementation still depends heavily on an accompanying programming language (in this case, C). This dependency prevents designers from creating a new Linda RPC operator for arbitrary applications of RPC. Therefore, every time a designer uses Linda for RPC, they must copy the source code for RPC or make a C-macro. This causes tight coupling, because the client must know the name of the RPC server. If the server name is passed in as a parameter, flexibility increases; however, this requires a binding phase in which the name is obtained and applied outside of the Linda framework.

The need remains for a design and programming methodology that allows implementation of communication protocols without tight coupling between the protocol implementation and the software elements with which the protocol implementation works.

A tuple space can require large quantities of dynamically allocated memory. However, most systems, and especially embedded systems, must operate within predictable and sometimes small memory requirements. Tuple-space systems are usually not suitable for coordination in systems that must operate within small predictable memory requirements because once a tuple has been generated, it remains in tuple space until it is explicitly removed or the software element that created it terminates. Maintaining a global tuple space can be very expensive in terms of overall system performance. Although much work has gone into improving the efficiency of tuple-space languages, system performance remains worse with tuple-space languages than with message-passing techniques.

The need remains for a design and programming methodology that can effectively coordinate between software elements while respecting performance and predictable memory requirements.

B. Fixed Coordination Models

5 In tuple-space languages, much of the complexity of coordination remains entangled with the functionality of computational elements. An encapsulating coordination formalism decouples inter-component interactions from the computational elements.

10 This type of formalism can be provided by fixed coordination models in which the coordination style is embodied in an entity and separated from computational concerns. Synchronous coordination models coordinate activity through relative schedules. Typically, these approaches require the coordination protocol to be manually constructed in advance. In addition, computational elements must be tailored to the coordination style used for a particular system (which may require
15 intrusive modification of the software elements).

The need remains for a design and programming methodology that allows for coordination between software elements without tailoring the software elements to the specific coordination style used in a particular software system while allowing for interactions between software elements in a way that facilitates debugging complex
20 systems.

Summary of the Invention

The present invention provides a coordination-centric debugging approach and programming methodology to facilitate the debugging of distributed software environments. This approach includes using "cooperative execution," in which the
25 software produces event traces to be analyzed by a debugging host.

A distributed software environment contains multiple processing elements, which are loaded with software, as well as sensors, displays, and other hardware. In accordance with the present invention, these processing elements contain software programs that generate corresponding event records in response to selected events. A

software program includes components, runtime systems, coordinators, interfaces, and runtime debugging architectures. The coordinator manages control and data flow interactions between components, and the interface between the coordinator and a component facilitates the exposure of an event. The runtime system of a processing
5 element collects the event record and transfers it to the runtime debugging architecture. The runtime debugging architecture facilitates the transfer of the event record from the processing element to a debugging host, along a communication channel. The communication channel facilitates either a direct or an indirect connection between the processing element and the debugging host.

10 In a direct connection, the runtime system of a processing element sends an event record to the runtime debugging architecture, which in turn connects to the debugging host and transfers the event record to the debugging host, which queues all event records and writes them to disk. In an indirect connection, an event record is routed from a primary runtime debugging architecture to an intermediate processing
15 element that collects the event record and transfers it to the debugging host.

When it is not possible or practical to connect directly or indirectly to a debugging host, event records may still be collected and debugged at a later time. A processing element may include a flash driver that interfaces with a flash memory. The processing element's runtime system sends the event records to a flash driver that
20 facilitates the collection and storage of event records in a flash memory for subsequent distributed debugging (*i.e.*, "post-mortem" debugging).

Each event record, whether for direct or indirect connections or for post-mortem debugging, is time-stamped and causality-stamped by the runtime debugging architecture before it is transferred to the debugging host, the intermediate processing
25 element, or flash memory. Time-stamping and causality-stamping facilitate coordination-centric debugging by stamping each event with a time in which the event was generated and with an identification of the cause of the event, allowing designers to debug a complex distributed software environment without having to guess at the proper event ordering.

In accordance with the present invention, any event made visible to the runtime system of a distributed software environment, either by being identified as an explicit event in the software or by event recording calls at runtime, is a candidate for recording. A distributed software environment contains instrumented distributed
5 code, which produces the events to be captured by the runtime system. The code includes information about how different software components interact. Instrumentation may consist of inserting event recording calls at each significant source line, which can cause a probe effect in the software but which effect can be minimized by selectively focusing the instrumentation. Also, events often occur in
10 specific, predetermined, partially ordered sequences; in some instances, a token representing the sequence may be recorded, rather than every single event in the sequence, thereby leaving the debugging host to expand the token.

In accordance with the present invention, the coordination-centric debugging approach makes complex distributed software environments more debuggable. The
15 debugging approach of the present invention facilitates the debugging of complex distributed systems and lets a designer isolate the cause, or causes, of unexpected system performance without having looked at, or modified, any of the underlying source code and avoid the time-consuming and error-prone process of source level debugging.

20 Additional aspects and advantages of this invention will be apparent from the following detailed description of preferred embodiments thereof, which proceeds with reference to the accompanying drawings.

Brief Description of the Drawings

Fig. 1 is a component in accordance with the present invention.

25 Fig. 2 is the component of Fig. 1 further having a set of coordination interfaces.

Fig. 3A is a prior art round-robin resource allocation protocol with a centralized controller.

Fig. 3B is a prior art round-robin resource allocation protocol implementing a token-passing scheme.

Fig. 4A is a detailed view of a component and a coordination interface connected to the component for use in round-robin resource allocation in accordance
5 with the present invention.

Fig. 4B depicts a round-robin coordinator in accordance with the present invention.

Fig. 5 shows several typical ports for use in a coordination interface in accordance with the present invention.

Fig. 6A is a unidirectional data transfer coordinator in accordance with the
10 present invention.

Fig. 6B is a bidirectional data transfer coordinator in accordance with the present invention.

Fig. 6C is a state unification coordinator in accordance with the present
15 invention.

Fig. 6D is a control state mutex coordinator in accordance with the present invention.

Fig. 7 is a system for implementing subsumption resource allocation having components, a shared resource, and a subsumption coordinator.

Fig. 8 is a barrier synchronization coordinator in accordance with the present
20 invention.

Fig. 9 is a rendezvous coordinator in accordance with the present invention.

Fig. 10 depicts a dedicated RPC system having a client, a server, and a dedicated RPC coordinator coordinating the activities of the client and the server.

Fig. 11 is a compound coordinator with both preemption and round-robin
25 coordination for controlling the access of a set of components to a shared resource.

Fig. 12A is software system with two data transfer coordinators, each having constant message consumption and generation rules and each connected to a separate data-generating component and connected to the same data-receiving component.

Fig. 12B is the software system of Fig. 12A in which the two data transfer
5 coordinators have been replaced with a merged data transfer coordinator.

Fig. 13 is a system implementing a first come, first served resource allocation protocol in accordance with the present invention.

Fig. 14 is a system implementing a multiclient RPC coordination protocol formed by combining the first come, first served protocol of Fig. 13 with the
10 dedicated RPC coordinator of Fig. 10.

Fig. 15 depicts a large system in which the coordination-centric design methodology can be employed having a wireless device interacting with a cellular network.

Fig. 16 shows a top-level view of the behavior and components of a system for
15 a cell phone.

Fig. 17A is a detailed view of a GUI component of the cell phone of Fig. 16.

Fig. 17B is a detailed view of a call log component of the cell phone of
Fig. 16.

Fig. 18A is a detailed view of a voice subsystem component of the cell phone
20 of Fig. 16.

Fig. 18B is a detailed view of a connection component of the cell phone of
Fig. 16.

Fig. 19 depicts the coordination layers between a wireless device and a base station, and between the base station and a switching center, of Fig. 15.

Fig. 20 depicts a cell phone call management component, a master switching
25 center call management component, and a call management coordinator connecting the respective call management components.

Fig. 21A is a detailed view of a transport component of the connection component of Fig. 18B.

Fig. 21B is a CDMA data modulator of the transport component of Fig. 18B.

Fig. 22 is a detailed view of a typical TDMA and a typical CDMA signal for the cell phone of Fig. 16.

Fig. 23A is a LCD touch screen component for a Web browser GUI for a wireless device.

Fig. 23B is a Web page formatter component for the Web browser GUI for the wireless device.

Fig. 24A is a completed GUI system for a handheld Web browser.

Fig. 24B shows the GUI system for the handheld Web browser combined with the connection subsystem of Fig. 18B in order to access the cellular network of Fig. 15.

Fig. 25 is a typical space/time diagram with space represented on a vertical axis and time represented on a horizontal axis.

Fig. 26 is a space/time diagram depicting a set of system events and two different observations of those system events.

Fig. 27 is a space/time diagram depicting a set of system events and an ideal observation of the events taken by a real-time observer.

Fig. 28 is a space/time diagram depicting two different yet valid observations of a system execution.

Fig. 29 is a space/time diagram depicting a system execution and an observation of that execution take by a discrete lamport observer.

Fig. 30 is a space/time diagram depicting a set of events that each include a lamport time stamp.

Fig. 31 is a space/time diagram illustrating the insufficiency of scalar timestamps to characterize causality between events.

Fig. 32 is a space/time diagram depicting a set of system events that each include a vector time stamp.

Fig. 33 depicts a display from a partial order event tracer (POET).

Fig. 34 is a space/time diagram depicting two compound events that are
5 neither causal nor concurrent.

Fig. 35 is a POET display of two convex event clusters.

Fig. 36 is a basis for distributed event environments (BEE) abstraction facility for a single client.

Fig. 37 is a hierarchical tree construction of process clusters.

10 Fig. 38A depicts a qualitative measure of cohesion and coupling between a set of process clusters that have heavy communication or are instantiated from the same source code.

Fig. 38B depicts a qualitative measure of cohesion and coupling between a set of process clusters that do not have heavy communication or are not instances of the
15 same source code.

Fig. 38C depicts a qualitative measure of cohesion and coupling between an alternative set of process clusters that have heavy communication or are instantiated from the same source code.

Fig. 39 depicts a consistent and an inconsistent cut of a system execution on a
20 space/time diagram.

Fig. 40A is a space/time diagram depicting a system execution.

Fig. 40B is a lattice representing all possible consistent cuts of the space/time diagram of Fig. 40A.

Fig. 40C is a graphical representation of the possible consistent cuts of Fig.
25 40B.

Fig. 41A is a space/time diagram depicting a system execution.

Fig. 41B is the space/time diagram of Fig. 41A after performing a global-step.

Fig. 41C is the space/time diagram of Fig. 41A after performing a step-over.

Fig. 41D is the space/time diagram of Fig. 41A after performing a step-in.

Fig. 42 is a space/time diagram depicting a system that is subject to a domino effect whenever the system is rolled back in time to a checkpoint.

5 Fig. 43 depicts the coordination-centric debugging approach in accordance with the present invention.

Fig. 44 is a detailed view of a direct connection between the primary processing element and the debugging host in accordance with the present invention.

10 Fig. 45 is a detailed view of an indirect connection between the primary processing element and the debugging host in accordance with the present invention.

Fig. 46 depicts capturing event records in flash memory for post-mortem distributed debugging in accordance with the present invention.

Fig. 47 shows how flash memory may be allocated.

15 Fig. 48 shows a distributed software environment being executed on a hardware platform and placement of a probe for monitoring bus traces on the platform and for generating event records.

Detailed Description of Preferred Embodiments

Coordination-Centric Software Design

20 Fig. 1 is an example of a component 100, which is the basic software element within the coordination-centric design framework, in accordance with the present invention. With reference to Fig. 1, component 100 contains a set of modes 102. Each mode 102 corresponds to a specific behavior associated with component 100. Each mode 102 can either be active or inactive, respectively enabling or disabling the behavior corresponding to that mode 102. Modes 102 can make the conditional
25 aspects of the behavior of component 100 explicit. The behavior of component 100 is encapsulated in a set of actions 104, which are discrete, event-triggered behavioral elements within the coordination-centric design methodology. Component 100 can be

copied and the copies of component 100 can be modified, providing the code-sharing benefits of inheritance.

Actions 104 are enabled and disabled by modes 102, and hence can be thought of as effectively being properties of modes 102. An event (not shown) is an instantaneous condition, such as a timer tick, a data departure or arrival, or a mode change. Actions 104 can activate and deactivate modes 102, thereby selecting the future behavior of component 100. This is similar to actor languages, in which methods are allowed to replace an object's behavior.

In coordination-centric design, however, all possible behaviors must be identified and encapsulated before runtime. For example, a designer building a user interface component for a cell phone might define one mode for looking up numbers in an address book (in which the user interface behavior is to display complete address book entries in formatted text) and another mode for displaying the status of the phone (in which the user interface behavior is to graphically display the signal power and the battery levels of the phone). The designer must define both the modes and the actions for the given behaviors well before the component can be executed.

Fig. 2 is component 100 further including a first coordination interface 200, a second coordination interface 202, and a third coordination interface 204.

Coordination-centric design's components 100 provide the code-sharing capability of object-oriented inheritance through copying. Another aspect of object-oriented inheritance is polymorphism through shared interfaces. In object-oriented languages, an object's interface is defined by its methods. Although coordination-centric design's actions 104 are similar to methods in object-oriented languages, they do not define the interface for component 100. Components interact through explicit and separate coordination interfaces, in this figure coordination interfaces 200, 202, and 204. The shape of coordination interfaces 200, 202, and 204 determines the ways in which component 100 may be connected within a software system. The way coordination interfaces 200, 202, and 204 are connected to modes 102 and actions 104 within component 100 determines how the behavior of component 100 can be

managed within a system. Systemwide behavior is managed through coordinators (see Fig. 4B and subsequent).

For our approach to be effective, several factors in the design of software elements must coincide: packaging, internal organization, and how elements coordinate their behavior. Although these are often treated as independent issues, conflicts among them can exacerbate debugging. We handle them in a unified framework that separates the internal activity from the external relationship of component 100. This lets designers build more modular components and encourages them to specify distributable versions of coordination protocols. Components can be reused in a variety of contexts, both distributed, and single processor 1.

1. Introduction to Coordination

Within this application, coordination refers to the predetermined ways by which components interact. Consider a common coordination activity: resource allocation. One simple protocol for this is round-robin: participants are lined up, and the resource is given to each participant in turn. After the last participant is served, the resource is given back to the first. There is a resource-scheduling period during which each participant gets the resource exactly once, whether or not it is needed.

Fig. 3A is prior art round-robin resource allocation protocol with a centralized controller 300, which keeps track of and distributes the shared resource (not shown) to each of software elements 302, 304, 306, 308, and 310 in turn. With reference to Fig. 3A, controller 300 alone determines which software element 302, 304, 306, 308, or 310 is currently allowed to use the resource and which has it next. This implementation of a round-robin protocol permits software elements 302, 304, 306, 308, and 310 to be modular, because only controller 300 keeps track of the software elements. Unfortunately, when this implementation is implemented on a distributed architecture (not shown), controller 300 must typically be placed on a single processing element (not shown). As a result, all coordination requests must go through that processing element, which can cause a communication performance bottleneck. For example, consider the situation in which software elements 304 and 306 are implemented on a first processing element (not shown) and controller 300 is

implemented on a second processing element. Software element 304 releases the shared resource and must send a message indicating this to controller 300. Controller 300 must then send a message to software element 306 to inform software element 306 that it now has the right to the shared resource. If the communication channel
5 between the first processing resource and the second processing resource is in use or the second processing element is busy, then the shared resource must remain idle, even though both the current resource holder and the next resource holder (software elements 304 and 306 respectively) are implemented on the first processing element (not shown). The shared resource must typically remain idle until communication can
10 take place and controller 300 can respond. This is an inefficient way to control access to a shared resource.

Fig. 3B is a prior art round-robin resource allocation protocol implementing a token passing scheme. With reference to Fig. 3B, this system consists of a shared resource 311 and a set of software elements 312, 314, 316, 318, 320, and 322. In
15 this system a logical token 324 symbolizes the right to access resource 311, *i.e.*, when a software element holds token 324, it has the right to access resource 311. When one of software elements 312, 314, 316, 318, 320, or 322 finishes with resource 311, it passes token 324, and with token 324 the access right, to a successor. This implementation can be distributed without a centralized controller, but as shown in
20 Figure v3B, this is less modular, because it requires each software element in the set to keep track of a successor.

Not only must software elements 312, 314, 316, 318, 320, and 322 keep track of successors, but each must implement a potentially complicated and error-prone protocol for transferring token 324 to its successor. Bugs can cause token 324 to be
25 lost or introduce multiple tokens 324. Since there is no formal connection between the physical system and complete topology maps (diagrams that show how each software element is connected to others within the system), some software elements might erroneously be serviced more than once per cycle, while others are completely neglected. However, these bugs can be extremely difficult to track after the system is
30 completed. The protocol is entangled with the functionality of each software element,

and it is difficult to separate the two for debugging purposes. Furthermore, if a few of the software elements are located on the same machine, performance of the implementation can be poor. The entangling of computation and coordination requires intrusive modification to optimize the system.

5 2. Coordination-Centric Design's Approach to Coordination

 The coordination-centric design methodology provides an encapsulating formalism for coordination. Components such as component 100 interact using coordination interfaces, such as first, second, and third coordination interfaces 200, 202, and 204, respectively. Coordination interfaces preserve component modularity
10 while exposing any parts of a component that participate in coordination. This technique of connecting components provides polymorphism in a similar fashion to subtyping in object-oriented languages.

 Fig. 4A is a detailed view of a component 400 and a resource access coordination interface 402 connected to component 400 for use in a round-robin coordination protocol in accordance with the present invention. With reference to
15 Fig. 4A, resource access coordination interface 402 facilitates implementation of a round-robin protocol that is similar to the token-passing round-robin protocol described above. Resource access coordination interface 402 has a single bit of control state, called access, which is shown as an arbitrated control port 404 that
20 indicates whether or not component 400 is holding a virtual token (not shown). Component 400 can only use a send message port 406 on access coordination interface 402 when arbitrated control port 404 is true. Access coordination interface 402 further has a receive message port 408.

 Fig. 4B show a round-robin coordinator 410 in accordance with the present
25 invention. With reference to Fig. 4B, round-robin coordinator 410 has a set of coordinator coordination interfaces 412 for connecting to a set of components 400. Each component 400 includes a resource access coordination interface 402. Each coordinator coordination interface 412 has a coordinator arbitrated control port 414, an incoming send message port 416 and an outgoing receive message port 418.
30 Coordinator coordination interface 412 in complimentary to resource access

coordination interface 402, and vice versa, because the ports on the two interfaces are compatible and can function to transfer information between the two interfaces.

The round-robin protocol requires round-robin coordinator 410 to manage the coordination topology. Round-robin coordinator 410 is an instance of more general
5 abstractions called coordination classes, in which coordination classes define specific coordination protocols and a coordinator is a specific implementation of the coordination class. Round-robin coordinator 410 contains all information about how components 400 are supposed to coordinate. Although round-robin coordinator 410 can have a distributed implementation, no component 400 is required to keep
10 references to any other component 400 (unlike the distributed round-robin implementation shown in Fig. 3B). All required references are maintained by round-robin coordinator 410 itself, and components 400 do not even need to know that they are coordinating through round-robin. Resource access coordination interface 402 can be used with any coordinator that provides the appropriate complementary interface.
15 A coordinator's design is independent of whether it is implemented on a distributed platform or on a monolithic single processor platform.

3. Coordination Interfaces

Coordination interfaces are used to connect components to coordinators. They are also the principle key to a variety of useful runtime debugging techniques.
20 Coordination interfaces support component modularity by exposing all parts of the component that participate in the coordination protocol. Ports are elements of coordination interfaces, as are guarantees and requirements, each of which will be described in turn.

A. Ports

25 A port is a primitive connection point for interconnecting components. Each port is a five-tuple (T; A; Q; D; R) in which:

- T represents the data type of the port. T can be one of int, boolean, char, byte, float, double, or cluster, in which cluster represents a cluster of data types (*e.g.*, an int followed by a float followed by two bytes).

- A is a boolean value that is true if the port is arbitrated and false otherwise.
- Q is an integer greater than zero that represents logical queue depth for a port.
- 5 • D is one of in, out, inout, or custom and represents the direction data flows with respect to the port.
- R is one of discard-on-read, discard-on-transfer, or hold and represents the policy for data removal on the port. Discard-on-read indicates that data is removed immediately after it is read (and any data in the logical queue are shifted), discard-on-transfer indicates that data is removed from
10 a port immediately after being transferred to another port, and hold indicates that data should be held until it is overwritten by another value. Hold is subject to arbitration.

Custom directionality allows designers to specify ports that accept or generate
15 only certain specific values. For example, a designer may want a port that allows other components to activate, but not deactivate, a mode. While many combinations of port attributes are possible, we normally encounter only a few. The three most common are message ports (output or input), state ports (output, input, or both; sometimes arbitrated), and control ports (a type of state port). Fig. 5 illustrates the
20 visual syntax used for several common ports throughout this application. With reference to Fig. 5, this figure depicts an exported state port 502, an imported state port 504, an arbitrated state port 506, an output data port 508, and an input data port 510.

1. Message Ports

25 Message ports (output and input) data ports 508 and 510 respectively) are either send (T; false; 1; out; discard-on-transfer) or receive (T; false; Q; in; discard-on-read). Their function is to transfer data between components. Data passed to a send port is transferred immediately to the corresponding receive port, thus it cannot be retrieved from the send port later. Receive data ports can have queues of

various depths. Data arrivals on these ports are frequently used to trigger and pass data parameters into actions. Values remain on receive ports until they are read.

2. State Ports

State ports take one of three forms:

- 5 1. (T; false; 1; out; hold)
2. (T; false; 1; in; hold)
3. (T; true; 1; inout; hold)

State ports, such as exported state port 502, imported state port 504, and arbitrated state port 506, hold persistent values, and the value assigned to a state port
10 may be arbitrated. This means that, unlike message ports, values remain on the state ports until changed. When multiple software elements simultaneously attempt to alter the value of arbitrated state port 506, the final value is determined based on arbitration rules provided by the designer through an arbitration coordinator (not shown).

15 State ports transfer variable values between scopes. In coordination-centric design, all variables referenced by a component are local to that component, and these variables must be explicitly declared in the component's scope. Variables can, however, be bound to state ports that are connected to other components. In this way, a variable value can be transferred between components and the variable value
20 achieves the system-level effect of a multivariable.

3. Control Ports

Control ports are similar to state ports, but a control port is limited to having the boolean data type. Control ports are typically bound to modes. Actions interact with a control port indirectly, by setting and responding to the values of a mode that
25 is bound to the control port.

For example, arbitrated control port 404 shown in Fig. 4A is a control port that can be bound to a mode (not shown) containing all actions that send data on a shared channel. When arbitrated control port 404 is false, the mode is inactive, disabling all actions that send data on the channel.

B. Guarantees

Guarantees are formal declarations of invariant properties of a coordination interface. There can be several types of guarantees, such as timing guarantees between events, guarantees between control state (*e.g.*, state A and state B are guaranteed to be mutually exclusive), etc. Although a coordination interface's guarantees reflect properties of the component to which the coordination interface is connected, the guarantees are not physically bound to any internal portions of the component. Guarantees can often be certified through static analysis of the software system. Guarantees are meant to cache various properties that are inherent in a component or a coordinator in order to simplify static analysis of the software system.

A guarantee is a promise provided by a coordination interface. The guarantee takes the form of a predicate promised to be invariant. In principle, guarantees can include any type of predicate (*e.g.*, $x > 3$, in which x is an integer valued state port, or $t_{ea} - t_{eb} < 2ms$). Throughout the remainder of this application, guarantees will be only event-ordering guarantees (guarantees that specify acceptable orders of events) or control-relationship guarantees (guarantees pertaining to acceptable relative component behaviors).

C. Requirements

A requirement is a formal declaration of the properties necessary for correct software system functionality. An example of a requirement is a required response time for a coordination interface—the number of messages that must have arrived at the coordination interface before the coordination interface can transmit, or fire, the messages. When two coordination interfaces are bound together, the requirements of the first coordination interface must be conservatively matched by the guarantees of the second coordination interface (*e.g.*, $x < 7$ as a guarantee conservatively matches $x < 8$ as a requirement). As with guarantees, requirements are not physically bound to anything within the component itself. Guarantees can often be verified to be sufficient for the correct operation of the software system in which the component is used. In sum, a requirement is a predicate on a first coordination interface that must

be conservatively matched with a guarantee on a complementary second coordination interface.

D. Conclusion Regarding Coordination Interfaces

A coordination interface is a four-tuple $(P; G; R; I)$ in which:

- 5 • P is a set of named ports.
- G is a set of named guarantees provided by the interface.
- R is a set of named requirements that must be matched by guarantees of connected interfaces.
- I is a set of named coordination interfaces.

- 10 As this definition shows, coordination interfaces are recursive. Coordinator coordination interface 412, shown in Fig. 4B, used for round-robin coordination is called *AccessInterface* and is defined in Table 1.

| Constituent | Value |
|--------------|--|
| ports | $P = \{ \text{access:StatePort, s:outMessagePort, r:inMessagePort} \}$ |
| guarantees | $G = \{ \neg \text{access} \Rightarrow \neg \text{s.gen} \}$ |
| requirements | $R = \emptyset$ |
| interfaces | $I = \emptyset$ |

- 15 Related to coordination interfaces is a recursive coordination interface descriptor, which is a five-tuple $(P_a; G_a; R_a; I_d; N_d)$ in which:
- P_a is a set of abstract ports, which are ports that may be incomplete in their attributes (*i.e.*, they do not yet have a datatype).
 - G_a is a set of abstract guarantees, which are guarantees between abstract ports.
 - 20 • R_a is a set of abstract requirements, which are requirements between abstract ports.
 - I_d is a set of coordination interface descriptors.

- N_d is an element of $Q \times Q$, where $Q = \{\infty\} \cup \mathbb{Z}^+$ and \mathbb{Z}^+ denotes the set of positive integers. N_d indicates the number or range of numbers of permissible interfaces.

Allowing coordination interfaces to contain other coordination interfaces is a powerful feature. It lets designers use common coordination interfaces as complex ports within other coordination interfaces. For example, the basic message ports described above are nonblocking, but we can build a blocking coordination interface (not shown) that serves as a blocking port by combining a wait state port with a message port.

10 4. Coordinators

A coordinator provides the concrete representations of intercomponent aspects of a coordination protocol. Coordinators allow a variety of static analysis debugging methodologies for software systems created with the coordination-centric design methodology. A coordinator contains a set of coordination interfaces and defines the relationships the coordination interfaces. The coordination interfaces complement the component coordination interfaces provided by components operating within the protocol. Through matched interface pairs, coordinators effectively describe connections between message ports, correlations between control states, and transactions between components.

For example, round-robin coordinator 410, shown in Fig. 4B, must ensure that only one component 400 has its component control port 404's value, or its access bit, set to true. Round-robin coordinator 410 must further ensure that the correct component 400 has its component control port 404 set to true for the chosen sequence. This section presents formal definitions of the parts that comprise coordinators: modes, actions, bindings, action triples, and constraints. These definitions culminate in a formal definition of coordinators.

A. Modes

A mode is a boolean value that can be used as a guard on an action. In a coordinator, the mode is most often bound to a control port in a coordination interface for the coordinator. For example, in round-robin coordinator 410, the modes of

concern are bound to a coordinator control port 414 of each coordinator coordination interface 412.

B. Actions

An action is a primitive behavioral element that can:

- 5
- Respond to events.
 - Generate events.
 - Change modes.

Actions can range in complexity from simple operations up to complicated pieces of source code. An action in a coordinator is called a transparent action
10 because the effects of the action can be precomputed and the internals of the action are completely exposed to the coordination-centric design tools.

C. Bindings

Bindings connect input ports to output ports, control ports to modes, state ports to variables, and message ports to events. Bindings are transparent and passive.
15 Bindings are simply conduits for event notification and data transfer. When used for event notification, bindings are called triggers.

D. Action Triples

To be executed, an action must be enabled by a mode and triggered by an event. The combination of a mode, trigger, and action is referred to as an action
20 triple, which is a triple (m; t; a) in which:

- m is a mode.
- t is a trigger.
- a is an action.

The trigger is a reference to an event type, but it can be used to pass data into
25 the action. Action triples are written: mode : trigger : action

A coordinator's actions are usually either pure control, in which both the trigger and action performed affect only control state, or pure data, in which both the trigger and action performed occur in the data domain. In the case of round-robin coordinator 410, the following set of actions is responsible for maintaining the
30 appropriate state:

$\text{access}_i : - \text{access}_i : + \text{access}_{(i+1) \bmod n}$

The symbol "+" signifies a mode's activation edge (*i.e.*, the event associated with the mode becoming true), and the symbol "-" signifies its deactivation edge.

When any coordinator coordination interface 412 deactivates its arbitrated control port 404's, access bit, the access bit of the next coordinator coordination interface 412 is automatically activated.

E. Constraints

In this dissertation, constraints are boolean relationships between control ports. They take the form:

10 Condition \Rightarrow Effect

This essentially means that the Condition (on the left side of the arrow) being true implies that Effect (on the right side of the arrow) is also true. In other words, if Condition is true, then Effect should also be true.

A constraint differs from a guarantee in that the guarantee is limited to communicating in-variant relationships between components without providing a way to enforce the in-variant relationship. The constraint, on the other hand, is a set of instructions to the runtime system dealing with how to enforce certain relationships between components. When a constraint is violated, two corrective actions are available to the system: (1) modify the values on the left-hand side to make the left-hand expression evaluate as false (an effect termed backpressure) or (2) alter the right-hand side to make it true. We refer to these techniques as LHM (left-hand modify) and RHM (right-hand modify). For example, given the constraint $x \Rightarrow \neg y$ and the value $x \wedge y$, with RHM semantics the runtime system must respond by disabling y or setting y to false. Thus, the value of $\neg y$ is set to true.

25 The decision of whether to use LHM, to use RHM, or even to suspend enforcement of a constraint in certain situations can dramatically affect the efficiency and predictability of the software system. Coordination-centric design does not attempt to solve simultaneous constraints at runtime. Rather, runtime algorithms use local ordered constraint solutions. This, however, can result in some constraints being violated and is discussed further below.

Round-robin coordinator 410 has a set of safety constraints to ensure that there is never more than one token in the system:

$$\text{access}_i \Rightarrow \forall_{j \neq i} \neg \text{access}_j$$

The above equation translates roughly as access_i implies not access_j for the set of all access_j where j is not equal to i . Even this simple constraint system can cause problems with local resolution semantics (as are LHM and RHM). If the runtime system attempted to fix all constraints simultaneously, all access modes would be shut down. If they were fixed one at a time, however, any duplicate tokens would be erased on the first pass, satisfying all other constraints and leaving a single token in the system.

Since high-level protocols can be built from combinations of lower-level protocols, coordinators can be hierarchically composed. A coordinator is a six-tuple $(I; M; B; N; A; X)$ in which:

- I is a set of coordination interfaces.
- M is a set of modes.
- B is a set of bindings between interface elements (*e.g.*, control ports and message ports) and internal elements (*e.g.*, modes and triggers).
- N is a set of constraints between interface elements.
- A is a set of action triples for the coordinator.
- X is a set of subcoordinators.

Figs. 6A, 6B, 6C, and 6D show a few simple coordinators highlighting the bindings and constraints of the respective coordinators. With reference to Fig. 6A, a unidirectional data transfer coordinator 600 transfers data in one direction between two components (not shown) by connecting incoming receive message port 408 to outgoing receive message port 418 with a binding 602. With reference to Fig. 6B, bidirectional data transfer coordinator 604 transfers data back and forth between two components (not shown) by connecting incoming receive message port 408 to outgoing receive message port 418 with binding 602 and connecting send message port 406 to incoming send message port 416 with a second binding 602.

Unidirectional data transfer coordinator 600 and bidirectional data transfer

coordinator 604 simply move data from one message port to another. Thus, each coordinator consists of bindings between corresponding ports on separate coordination interfaces.

With reference to Fig. 6C, state unification coordinator 606 ensures that a state port a 608 and a state port b 610 are always set to the same value. State unification coordinator 606 connects state port a 608 to state port b 610 with binding 602. With reference to Fig. 6D, control state mutex coordinator 612 has a first constraint 618 and a second constraint 620 as follows:

(1) $c \Rightarrow \neg d$ and

10 (2) $d \Rightarrow \neg c$.

Constraints 618 and 620 can be restated as follows:

(1) A state port c 614 having a true value implies that a state port d 616 has a false value, and

15 (2) State port d 616 having a true value implies that state port c 614 has a false value.

A coordinator has two types of coordination interfaces: up interfaces that connect the coordinator to a second coordinator, which is at a higher level of design hierarchy and down interfaces that connect the coordinator either to a component or to a third coordinator, which is at a lower level of design hierarchy. Down interfaces have names preceded with "~". Round-robin coordinator 410 has six down coordination interfaces (previously referred to as coordinator coordination interface 412), with constraints that make the turning off of any coordinator control port 414 (also referred to as access control port) turn on the coordinator control port 414 of the next coordinator coordination interface 412 in line. Table 2 presents all constituents of the round-robin coordinator.

| Constituent | Value |
|-------------------------|------------------------------------|
| coordination interfaces | $I = \text{AccessInterface}_{1-6}$ |

| | |
|-----------------|--|
| modes | $M = \text{access}_{1-6}$ |
| bindings | $B = \forall_{1 \leq i \leq 6} (\sim \text{AccessInterface}_i.\text{access}, \text{access}_i) \cup$ |
| constraints | $N = \forall_{1 \leq i \leq 6} (\forall_{(1 \leq j \leq 6) \wedge (i \neq j)} \text{access}_i \Rightarrow \neg \text{access}_j)$ |
| actions | $A = \forall_{1 \leq i \leq 6} \text{access}_i : -\text{access}_i : +\text{access}_{(i+1) \bmod 6}$ |
| subcoordinators | $X = \emptyset$ |

This tuple describes an implementation of a round-robin coordination protocol for a particular system with six components, as shown in round-robin coordinator 410. We use a coordination class to describe a general coordination protocol that may not have a fixed number of coordinator coordination interfaces. The coordination class is a six-tuple (Ic; Mc; Bc; Nc; Ac; Xc) in which:

- Ic is a set of coordination interface descriptors in which each descriptor provides a type of coordination interface and specifies the number of such interfaces allowed within the coordination class.
- Mc is a set of abstract modes that supplies appropriate modes when a coordination class is instantiated with a fixed number of coordinator coordination interfaces.
- Bc is a set of abstract bindings that forms appropriate bindings between elements when the coordination class is instantiated.
- Nc is a set of abstract constraints that ensures appropriate constraints between coordination interface elements are in place as specified at instantiation.
- Ac is a set of abstract action triples for the coordinator.
- Xc is a set of coordination classes (hierarchy).

While a coordinator describes coordination protocol for a particular application, it requires many aspects, such as the number of coordination interfaces and datatypes, to be fixed. Coordination classes describe protocols across many applications. The use of the coordination interface descriptors instead of coordination interfaces lets coordination classes keep the number of interfaces and datatypes

undetermined until a particular coordinator is instantiated. For example, a round-robin coordinator contains a fixed number of coordinator coordination interfaces with specific bindings and constraints between the message and state ports on the fixed number of coordinator coordination interfaces. A round-robin coordination class contains descriptors for the coordinator coordination interface type, without stating how many coordinator coordination interfaces, and instructions for building bindings and constraints between ports on the coordinator coordination interfaces when a particular round-robin coordinator is created.

5. Components

10 A component is a six-tuple (I; A; M; V; S; X) in which:

- I is a set of coordination interfaces.
- A is a set of action triples.
- M is a set of modes.
- V is a set of typed variables.
- 15 • S is a set of subcomponents.
- X is a set of coordinators used to connect the subcomponents to each other and to the coordination interfaces.

Actions within a coordinator are fairly regular, and hence a large number of actions can be described with a few simple expressions. However, actions within a component are frequently diverse and can require distinct definitions for each individual action. Typically, a component's action triples are represented with a table that has three columns: one for the mode, one for the trigger, and one for the action code. Table 3 shows some example actions from a component that can use round-robin coordination.

25

| Mode | Trigger | Action |
|---------------|---------|---|
| access | tick | AccessInterface.s.send("Test message"); -access; |
| \neg access | tick | waitCount ++; |

A component resembles a coordinator in several ways (for example, the modes and coordination interfaces in each are virtually the same). Components can have internal coordinators, and because of the internal coordinators, components do not always require either bindings or constraints. In the following subsections, various aspects of components are described in greater detail. These aspects of components include variable scope, action transparency, and execution semantics for systems of actions.

A. Variable Scope

To enhance a component's modularity, all variables accessed by an action within the component are either local to the action, local to the immediate parent component of the action, or accessed by the immediate parent component of the action via state ports in one of the parent component's coordination interfaces. For a component's variables to be available to a hierarchical child component, they must be exported by the component and then imported by the child of the component.

B. Action Transparency

An action within a component can be either a transparent action or an opaque action. Transparent and opaque actions each have different invocation semantics. The internal properties, *i.e.* control structures, variable, changes in state, operators, etc., of transparent actions are visible to all coordination-centric design tools. The design tools can separate, observe, and analyze all the internal properties of opaque actions. Opaque actions are source code. Opaque actions must be executed directly, and looking at the internal properties of opaque actions can be accomplished only through traditional, source-level debugging techniques. An opaque action must explicitly declare any mode changes and coordination interfaces that the opaque action may directly affect.

C. Action Execution

An action is triggered by an event, such as data arriving or departing a message port, or changes in value being applied to a state port. An action can change the value of a state port, generate an event, and provide a way for the software system

to interact with low-level device drivers. Since actions typically produce events, a single trigger can be propagated through a sequence of actions.

6. Protocols Implemented with Coordination Classes

In this section, we describe several coordinators that individually implement some common protocols: subsumption, barrier synchronization, rendezvous, and dedicated RPC.

A. Subsumption Protocol

A subsumption protocol is a priority-based, preemptive resource allocation protocol commonly used in building small, autonomous robots, in which the shared resource is the robot itself.

Fig. 7 shows a set of coordination interfaces and a coordinator for implementing the subsumption protocol. With reference to Fig. 7, a subsumption coordinator 700 has a set of subsumption coordinator coordination interfaces 702, which have a subsume arbitrated coordinator control port 704 and an incoming subsume message port 706. Each subsume component 708 has a subsume component coordination interface 710. Subsume component coordination interface 710 has a subsume arbitrated component control port 712 and an outgoing subsume message port 714. Subsumption coordinator 700 and each subsume component 708 are connected by their respective coordination interfaces, 702 and 710. Each subsumption coordinator coordination interface 702 in subsumption coordinator 700 is associated with a priority. Each subsume component 708 has a behavior that can be applied to a robot (not shown). At any time, any subsume component 708 can attempt to assert its behavior on the robot. The asserted behavior coming from the subsume component 708 connected to the subsumption coordinator coordination interface 702 with the highest priority is the asserted behavior that will actually be performed by the robot. Subsume components 708 need not know anything about other components in the system. In fact, each subsume component 708 is designed to perform independently of whether their asserted behavior is performed or ignored.

Subsumption coordinator 700 further has a slave coordinator coordination interface 716, which has an outgoing slave message port 718. Outgoing slave

message port 718 is connected to an incoming slave message port 720. Incoming slave message port 720 is part of a slave coordination interface 722, which is connected to a slave 730. When a subsume component 708 asserts a behavior and that component has the highest priority, subsumption coordinator 700 will control
 5 slave 730 (which typically controls the robot) based on the asserted behavior.

The following constraint describes the basis of the subsumption coordinator 700's behavior:

$$10 \quad \text{subsume}_p \Rightarrow \bigwedge_{i=1}^{p-1} \neg \text{subsume}_i$$

This means that if any subsume component 708 has a subsume arbitrated component control port 712 that has a value of true, then all lower-priority subsume
 15 arbitrated component control ports 712 are set to false. An important difference between round-robin and subsumption is that in round-robin, the resource access right is transferred only when surrendered. Therefore, round-robin coordination has cooperative release semantics. However, in subsumption coordination, a subsume component 708 tries to obtain the resource whenever it needs to and succeeds only
 20 when it has higher priority than any other subsume component 708 that needs the resource at the same time. A lower-priority subsume component 708 already using the resource must surrender the resource whenever a higher-priority subsume component 708 tries to access the resource. Subsumption coordination uses preemptive release semantics, whereby each subsume component 708 must always be
 25 prepared to relinquish the resource.

Table 4 presents the complete tuple for the subsumption coordinator.

| Constituent | Value |
|-------------------------|--|
| coordination interfaces | $I = (\text{Subsume}_{1-n}) \cup (\text{Output})$ |
| modes | $M = \text{subsume}_{1-n}$ |
| bindings | $B = \bigvee_{1 \leq i \leq n} (\text{Subsume}_i.\text{subsume}, \text{subsume}_i) \cup$ |
| constraints | $N = \bigvee_{1 \leq i \leq n} (\bigvee_{(1 \leq j \leq i)} \text{subsume}_i \Rightarrow \neg \text{subsume}_j)$ |
| actions | $A = \emptyset$ |
| subcoordinators | $X = \emptyset$ |

B. Barrier Synchronization Protocol

Other simple types of coordination that components might engage in enforce synchronization of activities. An example is barrier synchronization, in which each component reaches a synchronization point independently and waits. Fig. 8 depicts a barrier synchronization coordinator 800. With reference to Fig. 8, barrier synchronization coordinator 800 has a set of barrier synchronization coordination interfaces 802, each of which has a coordinator arbitrated state port 804, named wait. Coordinator arbitrated state port 804 is connected to a component arbitrated state port 806, which is part of a component coordination interface 808. Component coordination interface 808 is connected to a component 810. When all components 810 reach their respective synchronization points, they are all released from waiting. The actions for a barrier synchronization coordinator with n interfaces are:

$$\bigwedge_{0 \leq i < n} \text{wait}_i : : \bigvee_{0 \leq j < n} \neg \text{wait}_j$$

In other words, when all wait modes (not shown) become active, each one is released. The blank between the two colons indicates that the trigger event is the guard condition becoming true.

C. Rendezvous Protocol

A resource allocation protocol similar to barrier synchronization is called rendezvous. Fig. 9 depicts a rendezvous coordinator 900 in accordance with the present invention. With reference to Fig. 9, rendezvous coordinator 900 has a rendezvous coordination interface 902, which has a rendezvous arbitrated state port 904. A set of rendezvous components 906, each of which may perform different functions or have vastly different actions and modes, has a rendezvous component coordination interface 908, which includes a component arbitrated state port 910. Rendezvous components 906 connect to rendezvous coordinator 900 through their respective coordination interfaces, 908 and 902. Rendezvous coordinator 900 further has a rendezvous resource coordination interface 912, which has a rendezvous resource arbitrated state port 914, also called available. A resource 916 has a resource coordination interface 918, which has a resource arbitrated state port 920. Resource 916 is connected to rendezvous coordinator 900 by their complementary coordination interfaces, 918 and 912 respectively.

With rendezvous-style coordination, there are two types of participants: resource 916 and several resource users, here rendezvous components 916. When resource 916 is available, it activates its resource arbitrated state port 920, also referred to as its available control port. If there are any waiting rendezvous components 916, one will be matched with the resource; both participants are then released. This differs from subsumption and round-robin in that resource 916 plays an active role in the protocol by activating its available control port 920.

The actions for rendezvous coordinator 900 are:

$available_i \wedge wait_j : \quad : -available_i, -wait_j$

This could also be accompanied by other modes that indicate the status after the rendezvous. With rendezvous coordination, it is important that only one component at a time be released from wait mode.

D. Dedicated RPC Protocol

A coordination class that differs from those described above is dedicated RPC. Fig. 10 depicts a dedicated RPC system. With reference to Fig. 10, a dedicated RPC coordinator 1000 has an RPC server coordination interface 1002, which includes an RPC server imported state port 1004, an RPC server output message port 1006, and an RPC server input message port 1008. Dedicated RPC coordinator 1000 is connected to a server 1010. Server 1010 has a server coordination interface 1012, which has a server exported state port 1014, a server input data port 1016, and a server output data port 1018. Dedicated RPC coordinator 1000 is connected to server 1010 through their complementary coordination interfaces, 1002 and 1012 respectively. Dedicated RPC coordinator 1000 further has an RPC client coordination interface 1020, which includes an RPC client imported state port 1022, an RPC client input message port 1024, and an RPC client output message port 1026. Dedicated RPC coordinator 1000 is connected to a client 1028 by connecting RPC client coordination interface 1020 to a complementary client coordination interface 1030. Client coordination interface 1030 has a client exported state port 1032, a client output message port 1034, and a client input message port 1036.

The dedicated RPC protocol has a client/server protocol in which server 1010 is dedicated to a single client, in this case client 1028. Unlike the resource allocation protocol examples, the temporal behavior of this protocol is the most important factor in defining it. The following transaction listing describes this temporal behavior:

Client 1028 enters blocked mode by changing the value stored at client exported state port 1032 to true.

Client 1028 transmits an argument data message to server 1010 via client output message port 1034.

Server 1010 receives the argument (labeled "a") data message via server input data port 1016 and enters serving mode by changing the value stored in server exported state port 1014 to true.

Server 1010 computes return value.

Server 1010 transmits a return (labeled "r") message to client 1020 via server output data port 1018 and exits serving mode by changing the value stored in server exported state port 1014 to false.

Client 1028 receives the return data message via client input message port 1036 and exits blocked mode by changing the value stored at client exported state port 1032 to false.

This can be presented more concisely with an expression describing causal relationships:

$$\begin{aligned}
 10 \quad T_{RPC} = & +client.blocked \rightarrow client.transmits \rightarrow \\
 & +server.serving \rightarrow server.transmits \rightarrow \\
 & (-server.serving \parallel client.receives) \rightarrow -client.blocked
 \end{aligned}$$

The transactions above describe what is supposed to happen. Other properties of this protocol must be described with temporal logic predicates.

$$\begin{aligned}
 & server.serving \Rightarrow client.blocked \\
 & server.serving \Rightarrow F(server.r.output) \\
 & server.a.input \Rightarrow F(server.serving)
 \end{aligned}$$

20

The r in $server.r.output$ refers to the server output data port 1018, also labeled as the r event port on the server, and the a in $server.a.input$ refers to server input data port 1016, also labeled as the a port on the server (see Fig. 10).

Together, these predicates indicate that (1) it is an error for server 1010 to be in serving mode if client 1028 is not blocked; (2) after server 1010 enters serving mode, a response message is sent or else an error occurs; and (3) server 1010 receiving a message means that server 1010 must enter serving mode. Relationships between control state and data paths must also be considered, such as:

(client.a \Rightarrow client.blocked)

In other words, client 1028 must be in blocked mode whenever it sends an argument message.

- 5 The first predicate takes the same form as a constraint; however, since dedicated RPC coordinator 1000 only imports the client:blocked and server:serving modes (*i.e.*, through RPC client imported state port 1022 and RPC server imported state port 1004 respectively), dedicated RPC coordinator 1000 is not allowed to alter these values to comply. In fact, none of these predicates is explicitly enforced by a
- 10 runtime system. However, the last two can be used as requirements and guarantees for interface type-checking.

7. System-Level Execution

- Coordination-centric design methodology lets system specifications be executed directly, according to the semantics described above. When components and coordinators are composed into higher-order structures, however, it becomes essential
- 15 to consider hazards that can affect system behavior. Examples include conflicting constraints, in which local resolution semantics may either leave the system in an inconsistent state or make it cycle forever, and conflicting actions that undo one another's behavior. In the remainder of this section, the effect of composition issues
- 20 on system-level executions is explained.

A. System Control Configurations

- A configuration is the combined control state of a system—basically, the set of active modes at a point in time. In other words, a configuration in coordination-centric design is a bit vector containing one bit for each mode in the
- 25 system. The bit representing a control state is true when the control state is active and false when the control state is inactive. Configurations representing the complete system control state facilitate reasoning on system properties and enable several forms of static analysis of system behavior.

B. Action-Trigger Propagation

Triggers are formal parameters for events. As mentioned earlier, there are two types of triggers: (1) control triggers, invoked by control events such as mode change requests, and (2) data flow triggers, invoked by data events such as message arrivals or departures. Components and coordinators can both request mode changes (on the modes visible to them) and generate new messages (on the message ports visible to them). Using actions, these events can be propagated through the components and coordinators in the system, causing a cascade of data transmissions and mode change requests, some of which can cancel other requests. When the requests, and secondary requests implied by them, are all propagated through the system, any requests that have not been canceled are confirmed and made part of the system's new configuration.

Triggers can be immediately propagated through their respective actions or delayed by a scheduling step. Recall that component actions can be either transparent or opaque. Transparent actions typically propagate their triggers immediately, although it is not absolutely necessary that they do so. Opaque actions typically must always delay propagation.

1. Immediate Propagation

Some triggers must be immediately propagated through actions, but only on certain types of transparent actions. Immediate propagation can often involve static precomputation of the effect of changes, which means that certain actions may never actually be performed. For example, consider a system with a coordinator that has an action that activates mode A and a coordinator with an action that deactivates mode B whenever A is activated. Static analysis can be used to determine in advance that any event that activates A will also deactivate B; therefore, this effect can be executed immediately without actually propagating it through A.

2. Delayed Propagation

Trigger propagation through opaque actions must typically be delayed, since the system cannot look into opaque actions to precompute their results. Propagation may be delayed for other reasons, such as system efficiency. For example, immediate

propagation requires tight synchronization among software components. If functionality is spread among a number of architectural components, immediate propagation is impractical.

C. A Protocol Implemented with a Compound Coordinator

5 Multiple coordinators are typically needed in the design of a system. The multiple coordinators can be used together for a single, unified behavior. Unfortunately, one coordinator may interfere with another's behavior.

Fig. 11 shows a combined coordinator 1100 with both preemption and round-robin coordination for controlling access to a resource, as discussed above. With
10 reference to Fig. 11, components 1102, 1104, 1106, 1108, and 1110 primarily use round-robin coordination, and each includes a component coordination interface 1112, which has a component arbitrated control port 1114 and a component output message port 1116. However, when a preemptor component 1120 needs the resource, preemptor component 1120 is allowed to grab the resource immediately. Preemptor
15 component 1120 has a preemptor component coordination interface 1122. Preemptor component coordination interface 1122 has a preemptor arbitrated state port 1124, a preemptor output message port 1126, and a preemptor input message port 1128.

All component coordination interfaces 1112 and preemptor component coordination interface 1122 are connected to a complementary combined coordinator
20 coordination interface 1130, which has a coordinator arbitrated state port 1132, a coordinator input message port 1134, and a coordinator output message port 1136. Combined coordinator 1100 is a hierarchical coordinator and internally has a round-robin coordinator (not shown) and a preemption coordinator (not shown). Combined coordinator coordination interface 1130 is connected to a coordination
25 interface to round-robin 1138 and a coordination interface to preempt 1140. Coordinator arbitrated state port 1132 is bound to both a token arbitrated control port 1142, which is part of coordination interface to round-robin 1138, and to a preempt arbitrated control port 1144, which is part of coordination interface to preempt 1140. Coordinator input message port 1134 is bound to an interface to a round-robin output

message port 1146, and coordinator output message port 1136 is bound to an interface to round-robin input message port 1148.

Thus, preemption interferes with the normal round-robin ordering of access to the resource. After a preemption-based access, the resource moves to the component
5 that in round-robin-ordered access would be the successor to preemptor component 1120. If the resource is preempted too frequently, some components may starve.

D. Mixing Control and Data in Coordinators

Since triggers can be control-based, data-based, or both, and actions can produce both control and data events, control and dataflow aspects of a system are
10 coupled through actions. Through combinations of actions, designers can effectively employ modal data flow, in which relative schedules are switched on and off based on the system configuration.

Relative scheduling is a form of coordination. Recognizing this and understanding how it affects a design can allow a powerful class of optimizations.
15 Many data-centric systems (or subsystems) use conjunctive firing, which means that a component buffers messages until a firing rule is matched. When matching occurs, the component fires, consuming the messages in its buffer that caused it to fire and generating a message or messages of its own. Synchronous data flow systems are those in which all components have only firing rules with constant message
20 consumption and generation.

Fig. 12A shows a system in which a component N1 1200 is connected to a component N3 1202 by a data transfer coordinator 1204 and a component N2 1206 is connected to component N3 1202 by a second data transfer coordinator 1208. Component N3 1202 fires when it accumulates three messages on a port c 1210 and
25 two messages on a port d 1212. On firing, component N3 1202 produces two messages on a port o 1214. Coordination control state tracks the logical buffer depth for these components. This is shown with numbers representing the logical queue depth of each port in Fig. 12.

Fig. 12B shows the system of Fig. 12A in which data transfer coordinator
30 1204 and second data transfer coordinator 1208 have been merged to form a merged

data transfer coordinator 1216. Merging the coordinators in this example provides an efficient static schedule for component firing. Merged data transfer coordinator 1216 fires component N1 1200 three times and component N2 1206 twice. Merged data transfer coordinator 1216 then fires component N3 1202 twice (to consume all
5 messages produced by component N1 1200 and component N2 1206).

Message rates can vary based on mode. For example, a component may consume two messages each time it fires in one mode and four each time it fires in a second mode. For a component like this, it is often possible to merge schedules on a configuration basis, in which each configuration has static consumption and
10 production rates for all affected components.

E. Coordination Transformations

In specifying complete systems, designers must often specify not only the coordination between two objects, but also the intermediate mechanism they must use to implement this coordination. While this intermediate mechanism can be as simple
15 as shared memory, it can also be another coordinator; hence coordination may be, and often is, layered. For example, RPC coordination often sits on top of a TCP/IP stack or on an IrDA stack, in which each layer coordinates with peer layers on other processing elements using unique coordination protocols. Here, each layer provides certain capabilities to the layer directly above it, and the upper layer must be
20 implemented in terms of them.

In many cases, control and communication synthesis can be employed to automatically transform user-specified coordination to a selected set of standard protocols. Designers may have to manually produce transformations for nonstandard protocols.

F. Dynamic Behavior with Compound Coordinators

Even in statically bound systems, components may need to interact in a fashion that appears dynamic. For example, RPC-style coordination often has multiple clients for individual servers. Here, there is no apparent connection between client and server until one is forged for a transaction. After the connection is forged, however,
30 the coordination proceeds in the same fashion as dedicated RPC.

Our approach to this is to treat the RPC server as a shared resource, requiring resource allocation protocols to control access. However, none of the resource allocation protocols described thus far would work efficiently under these circumstances. In the following subsections, an appropriate protocol for treating the RPC as a shared resource will be described and how that protocol should be used as part of a complete multiclient RPC coordination class—one that uses the same RPC coordination interfaces described earlier—will be discussed.

1. First Come/First Serve protocol (FCFS)

Fig. 13 illustrates a first come/first serve (FCFS) resource allocation protocol, which is a protocol that allocates a shared resource to the requester that has waited longest. With reference to Fig. 13, a FCFS component interface 1300 for this protocol has a request control port 1302, an access control port 1304 and a component outgoing message port 1306. A FCFS coordinator 1308 for this protocol has a set of FCFS interfaces 1310 that are complementary to FCFS component interfaces 1300, having a FCFS coordinator request control port 1312, a FCFS coordinator access port 1314, and a FCFS coordinator input message port 1316. When a component 1318 needs to access a resource 1320, it asserts request control port 1302. When granted access, FCFS coordinator 1308 asserts the appropriate FCFS coordinator access port 1314, releasing FCFS coordinator request control port 1312.

To do this, FCFS coordinator 1308 uses a rendezvous coordinator and two round-robin coordinators. One round-robin coordinator maintains a list of empty slots in which a component may be enqueued, and the other round-robin coordinator maintains a list showing the next component to be granted access. When an FCFS coordinator request control port 1312 becomes active, FCFS coordinator 1308 begins a rendezvous access to a binder action. When activated, this action maps the appropriate component 1318 to a position in the round-robin queues. A separate action cycles through one of the queues and selects the next component to access the server. As much as possible, FCFS coordinator 1308 attempts to grant access to resource 1320 to the earliest component 1318 having requested resource 1320, with

concurrent requests determined based on the order in the rendezvous coordinator of the respective components 1318.

2. Multiclient RPC

Fig. 14 depicts a multiclient RPC coordinator 1400 formed by combining
5 FCFS coordinator 1308 with dedicated RPC coordinator 1000. With reference to Fig. 14, a set of clients 1402 have a set of client coordination interfaces 1030, as shown in Fig. 10. In addition, multiclient RPC coordinator 1400 has a set of RPC client coordination interfaces 1020, as shown in Fig. 10. For each RPC client coordination interface 1020, RPC client input message port 1024, of RPC client coordination
10 interface 1020, is bound to the component outgoing message port 1306 of FCFS coordinator 1308. Message transfer action 1403 serves to transfer messages between RPC client input message port 1024 and component outgoing message port 1306. For coordinating the actions of multiple clients 1402, multiclient RPC coordinator 1400 must negotiate accesses to a server 1404 and keep track of the values returned by
15 server 1404.

F. Monitor Modes and Continuations

Features such as blocking behavior and exceptions can be implemented in the coordination-centric design methodology with the aid of monitor modes. Monitor modes are modes that exclude all but a selected set of actions called continuations,
20 which are actions that continue a behavior started by another action.

1. Blocking Behavior

With blocking behavior, one action releases control while entering a monitor mode, and a continuation resumes execution after the anticipated response event. Monitor mode entry must be immediate (at least locally), so that no unexpected
25 actions can execute before they are blocked by such a mode.

Each monitor mode has a list of actions that cannot be executed when it is entered. The allowed (unlisted) actions are either irrelevant or are continuations of the action that caused entry into this mode. There are other conditions, as well. This mode requires an exception action if forced to exit. However, this exception action is
30 not executed if the monitor mode is turned off locally.

When components are distributed over a number of processing elements, it is not practical to assume complete synchronization of the control state. In fact, there are a number of synchronization options available as detailed in Chou, P “Control Composition and Synthesis of Distributed Real-Time Embedded Systems”, Ph.D. dissertation, University of Washington, 1998.

2. Exception Handling

Exception actions are a type of continuation. When in a monitor mode, exception actions respond to unexpected events or events that signal error conditions. For example, multiclient RPC coordinator 1400 can bind $\neg client.blocked$ to a monitor mode and set an exception action on $+server.serving$. This will signal an error whenever the server begins to work when the client is not blocked for a response.

8. A Complete System Example

Figure 15 depicts a large-scale example system under the coordination-centric design methodology. With reference to Fig. 15, the large scale system is a bimodal digital cellular network 1500. Network 1500 is for the most part a simplified version of a GSM (global system for mobile communications) cellular network. This example shows in greater detail how the parts of coordination-centric design work together and demonstrates a practical application of the methodology. Network 1500 has two different types of cells, a surface cell 1502 (also referred to as a base station 1502) and a satellite cell 1504. These cells are not only differentiated by physical position, but by the technologies they use to share network 1500. Satellite cells 1504 use a code division multiple access (CDMA) technology, and surface cells 1502 use a time division multiple access (TDMA) technology. Typically, there are seven frequency bands reserved for TDMA and one band reserved for CDMA. The goal is for as much communication as possible to be conducted through the smaller TDMA cells, here surface cells 1502, because power requirements for a CDMA cells, here satellite cell 1504, increase with the number of users in the CDMA cell. Mobile units 1506, or wireless devices, can move between surface cells 1502, requiring horizontal handoffs between surface cells 1502. Several surface cells 1502 are typically

connected to a switching center 1508. Switching center 1508 is typically connected to a telephone network or the Internet 1512. In addition to handoffs between surface cells 1502, the network must be able to hand off between switching centers 1508. When mobile units 1506 leave the TDMA region, they remain covered by satellite cells 1504 via vertical handoffs between cells. Since vertical handoffs require changing protocols as well as changing base stations and switching centers, they can be complicated in terms of control.

Numerous embedded systems comprise the overall system. For example, switching center 1508 and base stations, surface cells 1502, are required as part of the network infrastructure, but cellular phones, handheld Web browsers, and other mobile units 1506 may be supported for access through network 1500. This section concentrates on the software systems for two particular mobile units 1506: a simple digital cellular phone (shown in Fig. 16) and a handheld Web browser (shown in Fig. 24). These examples require a wide variety of coordinators and reusable components. Layered coordination is a feature in each system, because a function of many subsystems is to perform a layered protocol. Furthermore, this example displays how the hierarchically constructed components can be applied in a realistic system to help manage the complexity of the overall design.

To begin this discussion, we describe the cellular phone in detail, focusing on its functional components and the formalization of their interaction protocols. We then discuss the handheld Web browser in less detail but highlight the main ways in which its functionality and coordination differ from those of the cellular phone. In describing the cellular phone, we use a top-down approach to show how a coherent system organization is preserved, even at a high level. In describing the handheld Web browser, we use a bottom-up approach to illustrate component reuse and bottom-up design.

A. Cellular Phone

Fig. 16 shows a top-level coordination diagram of the behavior of a cell phone 1600. Rather than using a single coordinator that integrates the components under a

single protocol, we use several coordinators in concert. Interactions between coordinators occur mainly within the components to which they connect.

With reference to Fig. 16, cell phone 1600 supports digital encoding of voice streams. Before it can be used, it must be authenticated with a home master switching center (not shown). This authentication occurs through a registered master switch for each phone and an authentication number from the phone itself. There are various authentication statuses, such as full access, grey-listed, or blacklisted. For cell phone 1600, real-time performance is more important than reliability. A dropped packet is not retransmitted, and a late packet is dropped since its omission degrades the signal less than its late incorporation.

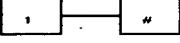
Each component of cell phone 1600 is hierarchical. A GUI 1602 lets users enter phone numbers while displaying them and query an address book 1604 and a logs component 1606. Address book 1604 is a database that can map names to phone numbers and vice versa. GUI 1602 uses address book 1604 to help identify callers and to look up phone numbers to be dialed. Logs 1606 track both incoming and outgoing calls as they are dialed. A voice component 1608 digitally encodes and decodes, and compresses and decompresses, an audio signal. A connection component 1610 multiplexes, transmits, receives, and demultiplexes the radio signal and separates out the voice stream and caller identification information.

Coordination among the above components makes use of several of the coordinators discussed above. Between connection component 1610 and a clock 1612, and between logs 1606 and connection component 1610, are unidirectional data transfer coordinators 600 as described with reference to Fig. 6A. Between voice component 1608 and connection component 1610, and between GUI 1602 and connection component 1610, are bidirectional data transfer coordinators 604, as described with reference to Fig. 6B. Between clock 1612 and GUI 1602 is a state unification coordinator 606, as described with reference to Fig. 6C. Between GUI 1602 and address book 1604 is a dedicated RPC coordinator 1000 as described with reference to Fig. 10, in which address book 1604 has client 1028 and GUI 1602 has server 1010.

There is also a custom GUI/log coordinator 1614 between logs 1606 and GUI 1602. GUI/log coordinator 1614 lets GUI 1602 transfer new logged information through an r output message port 1616 on a GUI coordination interface 1618 to an r input message port 1620 on a log coordination interface 1622. GUI/log coordinator 5 1614 also lets GUI 1602 choose current log entries through a pair of c output message ports 1624 on GUI coordination interface 1618 and a pair of c input message ports 1626 on log coordination interface 1622. Logs 1606 continuously display one entry each for incoming and outgoing calls.

1. GUI Component

10 Fig. 17A is a detailed view of GUI component 1602, of Fig. 16. With reference to Fig. 17A, GUI component 1602 has two inner components, a keypad 1700 and a text-based liquid crystal display 1702, as well as several functions of its own (not shown). Each time a key press occurs, it triggers an action that interprets the press, depending on the mode of the system. Numeric presses enter values into a 15 shared dialing buffer. When a complete number is entered, the contents of this buffer are used to establish a new connection through connection component 1610. Table 5 shows the action triples for GUI 1602.

| Mode | Trigger | Action |
|------------|---|---|
| Idle |  | numBuffer.append(keypress.val) |
| | Send | radio.send(numBuffer.val) + outgoingCall |
| | Disconnect | Nil |
| | Leftarrow | AddressBook.forward() + lookupMode |
| | =>Right | log.lastcall() + outlog |
| LookupMode | Leftarrow | AddressBook.forward() |
| | =>Right | AddressBook.backward() |

An "Addr Coord" coordinator 1704 includes an address book mode (not shown) in which arrow key presses are transformed into RPC calls.

2. Logs Component

5 Fig. 17B is a detailed view of logs component 1606, which tracks all incoming and outgoing calls. With reference to Fig. 17B, both GUI component 1602 and connection component 1610 must communicate with logs component 1606 through specific message ports. Those specific message ports include a transmitted number message port 1720, a received number message port 1722, a change current received message port 1724, a change current transmitted message port 1726, and two state
10 ports 1728 and 1729 for presenting the current received and current transmitted values, respectively.

Logs component 1606 contains two identical single-log components: a send log 1730 for outgoing calls and a receive log 1740 for incoming calls. The interface
15 of logs component 1606 is connected to the individual log components by a pair of adapter coordinators, Adap1 1750 and Adap2 1752. Adap1 1750 has an adapter receive interface 1754, which has a receive imported state port 1756 and a receive output message port 1758. Adap1 1750 further has an adapter send interface 1760,

which has a send imported state port 1762 and a send output message port 1764. Within Adap1, state port 1728 is bound to receive imported state port 1756, change current received message port 1724 is bound to receive output message port 1758, received number message port 1722 is bound to a received interface output message port 1766 on a received number coordination interface 1768, change current transmitted message port 1726 is bound to send output message port 1764, and state port 1729 is bound to Up.rc is bound to send imported state port 1762 .

3. Voice Component

Fig. 18A is a detailed view of voice component 1608 of Fig. 16. Voice component 1608 has a compression component 1800 for compressing digitized voice signals before transmission, a decompression component 1802 for decompressing received digitized voice signals, and interfaces 1804 and 1806 to analog transducers (not shown) for digitizing sound to be transmitted and for converting received transmissions into sound. Voice component 1608 is a pure data flow component containing sound generator 1808 which functions as a white-noise generator, a ring tone generator, and which has a separate port for each on sound generator interface 1810, and voice compression functionality in the form of compression component 1800 and decompression component 1802.

4. Connection Component

Fig. 18B is a detailed view of connection component 1610 of Fig. 16. With reference to Fig. 18B, connection component 1610 coordinates with voice component 1608, logs component 1606, clock 1612, and GUI 1602. In addition, connection component 1610 is responsible for coordinating the behavior of cell phone 1600 with a base station that owns the surface cell 1502 (shown in Fig. 15), a switching center 1508 (shown in Fig. 15), and all other phones (not shown) within surface cell 1502. Connection component 1610 must authenticate users, establish connections, and perform handoffs as needed—including appropriate changes in any low-level protocols (such as a switch from TDMA to CDMA).

Fig. 19 depicts a set of communication layers between connection component 1610 of cell phone 1600 and base station 1502 or switching center 1508. With

reference to Fig. 19, has several subcomponents, or lower-level components, each of which coordinates with an equivalent, or peer, layer on either base station 1502 or switching center 1508. The subcomponents of connection component 1610 include a cell phone call manager 1900, a cell phone mobility manager 1902, a cell phone radio resource manager 1904, a cell phone link protocol manager 1906, and a cell phone transport manager 1908 which is responsible for coordinating access to and transferring data through the shared airwaves TDMA and CDMA coordination. Each subcomponent will be described in detail including how each fits into the complete system.

10 Base station 1502 has a call management coordinator 1910, a mobility management coordinator 1912, a radio resource coordinator 1914 (BSSMAP 1915), a link protocol coordinator 1916 (SCCO 1917), and a transport coordinator 1918 (MTP 1919). Switching center 1508 has a switching center call manager 1920, a switching center mobility manager 1922, a BSSMAP 1924, a SCCP 1926, and an MTP 1928.

15 a. Call Management

Fig. 20 is a detailed view of a call management layer 2000 consisting of cell phone call manager 1900, which is connected to switching center call manager 1920 by call management coordinator 1910. With reference to Fig. 20, call management layer 2000 coordinates the connection between cell phone 1600 and switching center 1508. Call management layer 2000 is responsible for dialing, paging, and talking. Call management layer 2000 is always present in cell phone 1600, though not necessarily in Internet appliances (discussed later). Cell phone call manager 1900 includes a set of modes (not shown) for call management coordination that consists of the following modes:

- 25 • Standby
- Dialing
- RingingRemote
- Ringing
- CallInProgress

Cell phone call manager 1900 has a cell phone call manager interface 2002. Cell phone call manager interface 2002 has a port corresponding to each of the above modes. The standby mode is bound to a standby exported state port 2010. The dialing mode is bound to a dialing exported state port 2012. The RingingRemote mode is bound to a RingingRemote imported state port 2014. The Ringing mode is bound to a ringing imported state port 2016. The CallInProgress mode is bound to a CallInProgress arbitrated state port 2018.

Switching center call manager 1920 includes the following modes (not shown) for call management coordination at the switching center:

- 10 • Dialing
- RingingRemote
- Paging
- CallInProgress

Switching center call manager 1920 has a switching center call manager coordination interface 2040, which includes a port for each of the above modes within switching center call manager 1920.

When cell phone 1600 requests a connection, switching center 1508 creates a new switching center call manager and establishes a call management coordinator 1910 between cell phone 1600 and switching center call manager 1920.

20 b. Mobility Management

A mobility management layer authenticates mobile unit 1506 or cell phone 1600. When there is a surface cell 1502 available, mobility manager 1902 contacts the switching center 1508 for surface cell 1502 and transfers a mobile unit identifier (not shown) for mobile unit 1506 to switching center 1508. Switching center 1508 then looks up a home motor switching center for mobile unit 1506 and establishes a set of permissions assigned to mobile unit 1506. This layer also acts as a conduit for the call management layer. In addition, the mobility management layer performs handoffs between base stations 1502 and switching centers 1508 based on information received from the radio resource layer.

c. Radio Resource

In the radio resource layer, radio resource manager 1904, chooses the target base station 1502 and tracks changes in frequencies, time slices, and CDMA codes. Cell phones may negotiate with up to 16 base stations simultaneously. This layer also
5 identifies when handoffs are necessary.

d. Link Protocol

The link layer manages a connection between cell phone 1600 and base station 1502. In this layer, link protocol manager 1906 packages data for transfer to base station 1502 from cell phone 1600.

10 e. Transport

Fig. 21A is a detailed view of transport component 1908 of connection component 1610. Transport component 1908 has two subcomponents, a receive component 2100 for receiving data and a transmit component 2102 for transmitting data. Each of these subcomponents has two parallel data paths a CDMA path 2104
15 and a TDMA/FDMA path 2106 for communicating in the respective network protocols.

Fig. 21B is a detailed view of a CDMA modulator 2150, which implements a synchronous data flow data path. CDMA modulator 2150 takes the dot-product of an incoming data signal along path 2152 and a stored modulation code for cell phone
20 1600 along path 2154, which is a sequence of chips, which are measured time signals having a value of -1 or $+1$.

Transport component 1908 uses CDMA and TDMA technologies to coordinate access to a resource shared among several cell phones 1600, *i.e.*, the airwaves. Transport components 1908 supersede the FDMA technologies (*e.g.*, AM and FM)
25 used for analog cellular phones and for radio and television broadcasts. In FDMA, a signal is encoded for transmission by modulating it with a carrier frequency. A signal is decoded by demodulation after being passed through a band pass filter to remove other carrier frequencies. Each base station 1502 has a set of frequencies—chosen to minimize interference between adjacent cells. (The area covered by a cell may be
30 much smaller than the net range of the transmitters within it.)

TDMA, on the other hand, coordinates access to the airwaves through time slicing. Cell phone 1600 on the network is assigned a small time slice, during which it has exclusive access to the media. Outside of the small time slice, cell phone 1600 must remain silent. Decoding is performed by filtering out all signals outside of the small time slice. The control for this access must be distributed. As such, each component involved must be synchronized to observe the start and end of the small time slice at the same instant.

Most TDMA systems also employ FDMA, so that instead of sharing a single frequency channel, cell phones 1600 share several channels. The band allocated to TDMA is broken into frequency channels, each with a carrier frequency and a reasonable separation between channels. Thus user channels for the most common implementations of TDMA can be represented as a two-dimensional array, in which the rows represent frequency channels and the columns represent time slices.

CDMA is based on vector arithmetic. In a sense, CDMA performs inter-cell-phone coordination using data flow. Instead of breaking up the band into frequency channels and time slicing these, CDMA regards the entire band as an n-dimensional vector space. Each channel is a code that represents a basis vector in this space. Bits in the signal are represented as either 1 or -1, and the modulation is the inner product of this signal and a basis vector of mobile unit 1506 or cell phone 1600. This process is called spreading, since it effectively takes a narrowband signal and converts it into a broadband signal.

Demultiplexing is simply a matter of taking the dot-product of the received signal with the appropriate basis vector, obtaining the original 1 or -1. With fast computation and the appropriate codes or basis vectors, the signal can be modulated without a carrier frequency. If this is not the case, a carrier and analog techniques can be used to fill in where computation fails. If a carrier is used, however, all units use the same carrier in all cells.

Fig. 22 shows TDMA and CDMA signals for four cell phones 1600. With reference to Fig. 22, for TDMA, each cell phone 1600 is assigned a time slice during which it can transmit. Cell phone 1 is assigned time slice t_0 , cell phone 2 is assigned

time slice t1, cell phone 3 is assigned time slice t2, and cell phone 4 is assigned time slice t3. For CDMA, each cell phone 1600 is assigned a basis vector that it multiplies with its signal. Cell phone 1 is assigned the vector:

5

$$\begin{pmatrix} -1 \\ 1 \\ -1 \\ 1 \end{pmatrix}$$

Cell phone 2 is assigned the vector:

$$\begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix}$$

10

Cell phone 3 is assigned the vector:

$$\begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix}$$

Cell phone 4 is assigned the vector:

$$\begin{pmatrix} -1 \\ -1 \\ 1 \\ 1 \end{pmatrix}$$

5

Notice that these vectors form an orthogonal basis.

B. Handheld Web Browser

In the previous subsection, we demonstrated our methodology on a cell phone with a top-down design approach. In this subsection, we demonstrate our methodology with a bottom-up approach in building a handheld Web browser.

Fig. 23A is a LCD touch screen component 2300 for a Web browser GUI (shown in Fig. 24A) for a wireless device 1506. With reference to Fig. 23A, a LCD touch screen component 2300, has an LCD screen 2302 and a touch pad 2304.

Fig. 23B is a Web page access component 2350 for fetching and formatting web pages. With reference to Fig. 23B, web access component 2350 has a page fetch subcomponent 2352 and a page format subcomponent 2354. Web access component 2350 reads hypertext markup language (HTML) from a connection interface 2356, sends word placement requests to a display interface 2358, and sends image requests to the connection interface 2356. Web access component 2350 also has a character input interface to allow users to enter page requests directly and to fill out forms on pages that have forms.

Fig. 24A shows a completed handheld Web browser GUI 2400. With reference to Fig. 24A, handheld Web browser GUI 2400, has LCD touch screen component 2300, web access component 2350, and a pen stroke recognition component 2402 that translates pen strokes entered on touch pad 2304 into characters.

Fig. 24B shows the complete component view of a handheld Web browser 2450. With reference to Fig. 24B, handheld Web browser 2450 is formed by connecting handheld Web browser GUI 2400 to connection component 1610 of cell phone 1600 (described with reference to Fig. 16) with bi-directional data transfer coordinator 604 (described with reference to Fig. 6B). Handheld Web browser 2450 is an example of mobile unit 1506, and connects to the Internet through the cellular infrastructure described above. However, handheld Web browser 2450 has different access requirements than does cell phone 1600. For handheld Web browser 2450, reliability is more important than real-time delivery. Dropped packets usually require retransmission, so it is better to deliver a packet late than to drop it. Real-time issues primarily affect download time and are therefore secondary. Despite this, handheld Web browser 2450 must coordinate media access with cell phones 1600, and so it must use the same protocol as cell phones 1600 to connect to the network. For that reason, handheld Web browser 2450 can reuse connection component 1610 from cell phone 1600.

Debugging Techniques

In concept, debugging is a simple process. A designer locates the cause of undesired behavior in a system and fixes the cause. In practice, debugging—even of sequential software—remains difficult. Embedded systems are considerably more complicated to debug than sequential software, due to factors such as concurrence, distributed architectures, and real-time concerns. Issues taken for granted in sequential software, like a schedule that determines the order of all events (the program), are nonexistent in a typical distributed system. Locating and fixing bugs in these complex systems requires many factors including an understanding of the thought processes underpinning the design.

Prior art research into debugging distributed systems is diverse, eclectic, and lacks any standard notations. This application uses a standardized notation both to describe the prior art and the present invention. As a result of this standardized notation, the principles in the prior art follow those published in the referenced works. However, the specific notation, theorems, etc. may differ.

The two general classes of debugging techniques are event based debugging and state based debugging. Most debugging techniques for general-purpose distributed systems are event based. Event based debugging techniques operate by collecting event traces from individual system components and causally relating those event traces. These techniques require an ability to determine efficiently the causal ordering among any given pair of events. Determining the causal order can be difficult and costly.

Events may be primitive, or they may be hierarchical clusters of other events. Primitive events are abstractions for individual local occurrences that might be important to a debugger. Examples of primitive events in sequential programs are variable assignments and subroutine entry or returns. Primitive events for distributed systems include message send and receive events.

State-based debugging techniques are less commonly used in debugging distributed systems. State-based debugging techniques typically operate by presenting designers with views or snapshots of process state. Distributed systems are not tightly synchronized, and so these techniques traditionally involve only the state of individual processes. However, state-based debugging techniques can be applied more generally by relaxing the concept of an "instant in time" so that it can be effectively applied to asynchronous processes.

1. Event-based debugging

In this section, prior art systems for finding and tracking meaningful event orderings, despite limits in observation, are described. Typical ways in which event orderings are used in visualization tools through automated space/time diagrams are then described.

A. Event order determination and observation

The behavior of a software system is determined by the events that occur and the order in which they occur. For sequential systems, this seems almost too trivial to mention; of course, a given set of events, such as

$$\{x := 2, x := x * 2, x := 5, y := x\},$$

arranged in two different ways may describe two completely different behaviors.

However, since a sequential program is essentially a complete schedule of events,

- 5 ordering is explicit. Sequential debugging tools depend on the invariance of this event schedule to let programmers reproduce failures by simply using the same inputs. In distributed systems, as in any concurrent system, it is neither practical nor efficient to completely schedule all events. Concurrent systems typically must be designed with flexible event ordering.

- 10 Determining the order in which events occur in a distributed system is subject to the limits of observation. An observation is an event record collected by an observer. An observer is an entity that watches the progress of an execution and records events, but does not interfere with the system. To determine the order in which two events occur, an observer must measure them both against a common
15 reference.

- Fig. 25 shows a typical space/time diagram 2500, with space represented on a vertical axis 2502 and time represented on a horizontal axis 2504. With reference to Fig. 25, space/time diagram 2500 provides a starting point for discussing executions in distributed systems. Space/time diagram 2500 gives us a visual representation for
20 discussing event ordering and for comparing various styles of observation. A set of horizontal world lines 2506, 2508, and 2510 each represent an entity that is stationary in space. The entities represented by horizontal world lines 2506, 2508, and 2510 are called processes and typically represent software processes in the subject system. The entities can also represent any entity that generates events in a sequential fashion.
- 25 The spatial separation in the diagram, along vertical axis 2502, represents a virtual space, since several processes might execute on the same physical hardware. A diagonal world line 2512 is called a message and represents discrete communications that pass between two processes. A sphere 2514 represents an event. In subsequent figures vertical axis 2502 and horizontal axis 2504 are omitted from any space/time

diagrams, unless vertical axis 2502 and horizontal axis 2504 provide additional clarity to a particular figure.

Fig. 26 shows a space/time diagram 2600 of two different observations of a single system execution, taken by a first observer 2602 and a second observer 2604.

- 5 With reference to Fig. 26, first observer 2602 and second observer 2604 are entities that record event occurrence. First observer 2602 and second observer 2604 must each receive distinct notifications of each event that occurs and each must record the events in some total order. First observer 2602 and second observer 2604 are represented in space/time diagram 2600 as additional processes, or horizontal world
- 10 lines. Each event recorded requires a signal from its respective process to both first observer 2602 and second observer 2604. The signals from an event x 2606 on a process 2608 to both first observer 2602 and second observer 2604 are embodied in messages 2610 and 2612, respectively. First observer 2602 records event x 2606 as preceding an event y 2614. However, second observer 2604 records event y 2614 as
- 15 preceding event x 2606. Such effects may be caused by nonuniform latencies within the system.

- However, the observations of first observer 2602 and second observer 2604 are not equally valid. A valid observation is typically an observation that preserves the order of events that depend on each other. Second observer 2604 records the
- 20 receipt of a message 2616 before that message is transmitted. Thus the observation from second observer 2604 is not valid.

- Fig. 27 shows a space/time diagram 2700 for a special, ideal observer, called the real-time observer (RTO) 2702. With reference to Fig. 27, RTO 2702 can view each event immediately as it occurs. Due to the limitations of physical clocks, and
- 25 efficiency issues in employing them, it is usually not practical to implement RTO 2702. However, RTO 2702 represents an upper bound on precision in event-order determination.

Fig. 28 shows a space/time graph 2800 showing two valid observations of a system taken by two separate observers: RTO 2702 and a third observer 2802. With

reference to Fig. 28, there is nothing special about the ordering of the observation taken by RTO 2702. Events d 2804, e 2806, and f 2808 are all independent events in this execution. Therefore, the observation produced by RTO 2702 and the observation produced by third observer 2802 can each be used to reproduce equivalent
 5 executions of the system. Any observation in which event dependencies are preserved is typically equal in value to an observation by RTO 2702. However, real-time distributed systems may need additional processes to emulate timing constraints.

Fig. 29, is a space/time diagram 2900 of a methodological observer, called the discrete Lamport Observer (DLO) 2902, that records each event in a set of ordered
 10 bins. With reference to Fig. 29, DLO 2902 records an event 2904 in an ordered bin 2906 based on the following rule: each event is recorded in the leftmost bin that follows all events on which it depends. DLO 2902 views events discretely and does not need a clock. DLO 2902 does however require explicit knowledge of event dependency. To determine the bin in which each event must be placed, DLO 2902
 15 needs to know the bins of the immediately preceding events. The observation produced by DLO 2902 is also referred to as a topological sort of the system execution's event graph.

In the following, E is the set of all events in an execution. The immediate predecessor relation, $\sim \subseteq E \times E$, includes all pairs (e_a, e_b) such that:

20 a) If e_a and e_b are on the same process, e_a , precedes e_b with no intermediate events.

b) If e_b is a receive event, e_a is the send event that generated the message.

Given these conditions, e_a is called the immediate predecessor of e_b .

Each event has at most two immediate predecessors. Therefore, DLO 2902
 25 need only find the bins of at most two records before each placement. The transitive closure of the immediate predecessor relation forms a causal relation. The causal relation, $\sim^+ \subseteq E \times E$; is the smallest transitive relation such that $e_i \rightarrow e_j \Rightarrow e_i \sim^+ e_j$.

This relation defines a partial order of events and further limits the definition of a valid observation. A valid observation is an ordered record of events from a given execution, i.e., (R, \prec) , where $e \in E \Rightarrow (\text{record}(e)) \in R$, and \prec is an ordering operator. A valid observation has:

5

$$e_i, e_j \in E, e_i \sim e_j \Rightarrow \text{record}(e_i) \prec \text{record}(e_j)$$

The dual of the causal relation is a concurrence relation. The concurrence relation, $E \times E$, includes all pairs (e_a, e_b) such that neither $e_a \sim e_b$ nor $e_b \sim e_a$. While
10 the causal relation is transitive, the concurrence relation is not. The concurrence relation is symmetric, where the causal relation is not.

B. Event order tracking

Debugging typically requires an understanding of the order in which events
15 occur. Above, observers were presented as separate processes, while that treatment simplified the discussion of observers it is typically not a practical implementation of an observer. When the observer is implemented as a physical process the signals to indicate events would have to be transformed into physical messages, and the system would have to be synchronized to enable all messages to arrive in a valid order.

20 Fig. 30, depicts a space/time graph 3000 with each event having a label 3002. With reference to Fig. 30, DLO 2902 can accurately place event records in their proper bins—even if received out of order—as long as it knows the bins of the immediate predecessors. If we know the bins in which events are recorded, we can determine something about their causality. Fortunately, it is easy to label each event
25 with the number of its intended bin. Labels 3002 are analogous to time and are typically called Lamport timestamps.

A Lamport timestamp is an integer t associated with an event e_i such that

$$e_i \sim e_j \Rightarrow t(e_i) < t(e_j)$$

- 5 Lamport timestamps can be assigned as needed, provided the labels of an event's immediate predecessors are known. This information can be maintained with a local counter, called a Lamport clock (not shown), t_{Pi} , on each process, P_i . The clock's value is transmitted with each message M_j as t_{Mj} . Clock value t_{Pi} is updated with each event, as follows:

10

$$t_{pi} = \begin{cases} \max(t_{Mj}, t_{pi}) + 1; & \text{if } e \text{ is a receive event} \\ t_{pi} + 1 & ; \text{otherwise} \end{cases}$$

A labeling mechanism is said to characterize the causal relation if, based on their labels alone, it can be determined whether two events are causal or concurrent.

- 15 Although Lamport timestamps are consistent with causality (if $t(e_i) \geq t(e_j)$, then $e_i \rightarrow e_j$, they do not characterize the causal relation.

Fig. 31 is a space/ time graph 3100 that demonstrates the inability of scalar time stamps to characterize causality between events. With reference to Fig. 31, space/ time graph 3100 shows event e_1 3102, event e_2 3104, and event e_3 3106.

- 20 e_1 3102 \sim e_2 3104, and also $(e_1$ 3102 \parallel e_3 3106) \wedge (e_2 3104 \parallel e_3 3106) with scalar time stamps, $\text{label}(e_i)$, that characterize causality. It can be shown that e_3 3106 appears, when scalar timestamps are used, concurrent with both e_1 3102 and e_2 3104. However, since e_1 3102 \sim e_2 3104 it is not possible for e_3 3106 to be concurrent with both.

Event causality can be tracked completely using explicit event dependence graphs, with directed edges from each event to its immediate predecessors.

Unfortunately, this method cannot store enough information with each record to determine whether two arbitrarily chosen events are causally related without
5 traversing the dependence graph.

Other labeling techniques, such as vector timestamps, can characterize causality. The typical formulation of vector timestamps is based on the cardinality of event histories. A basis for vector timestamp is established based on the following definitions and theorems. An event history, $H(e_j)$, of an event e_j is the set of all

10 events, e_i , such that either since $e_i \rightsquigarrow e_j$ or $e_i \rightsquigarrow e_i = e_j$. The event history can be projected against specific processes. For a process P_j : the P_j history projection of $H(e_i)$, $H_{P_j}(e_i)$, is the intersection of $H(e_i)$ and the set of events local to P_j . The event graph represented by a space/time diagram can be partitioned into equivalence classes, with one class for each process. The set of events local to P_j is just the P_j
15 equivalence class.

The intersection of any two projections from the same process is identical to at least one of the two projections. Two history projections from a single process, $H_p(a)$ and $H_p(b)$, must satisfy one of the following:

- 20 a) $H_p(a) \subset H_p(b)$
 b) $H_p(a) = H_p(b)$
 c) $H_p(a) \supset H_p(b)$

25 The cardinality of $H_{P_j}(e_i)$ is thus the number of events local to P_j that causally precede e_i , and e_i itself. Since local events always occur in sequence, we can uniquely identify an event by its process and the cardinality of its local history.

For events $e_a; e_b$ with $e_a \neq e_b$, $H_{Pea}(e_a) \subseteq H_{Pea}(e_b) \Rightarrow e_a \sim e_b$

Fig. 32 shows a space/time diagram 3200 with vector timestamped events. A vector timestamp 3202 is a vector label, t_e , assigned to each event, $e \in E$, such that the i^{th} element represents $[H_{Pi}(e)]$. Given two events, e_1 and e_2 , we can determine their causal ordering: if vector t_{e_i} has a smaller value for its own process's entry than the other, t_{e_j} , has at that same position, then $e_i \sim e_j$. If both vectors have larger values for their own process entries, then $e_i \parallel e_j$. It is not possible for both events to have smaller values for their own entries because for events e_a and e_b , $e_a \rightarrow e_b$ implies $H_{Pea}(e_a) \supseteq H_{Pea}(e_b)$. It is not necessary to know the local processes of events to determine their causal order using vector timestamps.

The causal order of two vector timestamped events, e_a and e_b , from unknown processes can be determined with an element-by-element comparison of their vector timestamps:

$$\bigwedge_{i=1}^n t_{ea}[i] \leq t_{eb}[i] \Rightarrow e_a \rightarrow e_b$$

$$\neg \bigwedge_{i=1}^n t_{ea}[i] \leq t_{eb}[i] \quad \wedge$$

$$\neg \bigwedge_{i=1}^n t_{eb}[i] \leq t_{ea}[i] \Rightarrow e_a \parallel e_b$$

Thus, vector timestamps both fully characterize causality and uniquely identify each event in an execution.

Computing vector timestamps at runtime is similar to Lamport timestamp computation. Each process (P_s) contains a vector clock (\hat{t}_{Ps}) with elements for every

process in the system, where $\hat{t}_{Ps}[s]$ always equals the number of events local to P_s . Snapshots of this vector counter are used to label each event, and snapshots are transmitted with each message. The recipient of a message with a vector snapshot can update its own vector counter (\hat{t}_{Pr}) by replacing it with $\sup(\hat{t}_{Ps}, \hat{t}_{Pr})$, the element-wise maximum of \hat{t}_{Ps} and \hat{t}_{Pr} .

This technique places enough information with each message to determine message ordering. It is performed by comparing snapshots attached to each message. However, transmission of entire snapshots is usually not practical, especially if the system contains a large number of processes.

Vector clocks can however be maintained without transmitting complete snapshots. A transmitting process, P_s , can send a list that includes only those vector clock values that have changed since its last message. A recipient, P_r , then compares the change list to its current elements and updates those that are smaller. This requires each process to maintain several vectors: one for itself, and one for each process to which it has sent messages. However, change lists do not contain enough information to independently track message order.

The expense of maintaining vector clocks can be a strong deterrent to employing them. Unfortunately, no technique with smaller labels can characterize causality. It has been proven that the dimension of the causal relation for an N -process distributed execution is N , and hence N -element vectors are the smallest labels characterizing causality.

The problem results from concurrence, without which Lamport time would be sufficient. Concurrence can be tracked with concurrency maps, where each event keeps track of all events with which it is concurrent. Since the maps characterize concurrency, adding Lamport time lets them also characterize causality (the concurrency information disambiguates the scalar time). Unfortunately, concurrency maps can only be constructed after-the-fact, since doing so requires examination of events from all processes

In some situations, distinguishing between concurrency and causality is not a necessity, but merely a convenience. There are compact labeling techniques that allow better concurrence detection than Lamport time. One such technique uses interval clocks, in which each event record is labeled with its own Lamport time and the Lamport time of its earliest successor. This label then represents a Lamport time interval, during which the corresponding event was the latest known by the process. This gives each event a wider region with which to detect concurrence (indicated by overlapping intervals).

In cases where there is little or no cross-process causality (few messages), interval timestamps are not much better than Lamport timestamps. In cases with large numbers of messages, however, interval timestamps can yield better results.

C. Space/time displays in debugging tools

Space/time diagrams have typically proven useful in discussing event causality and concurrence. Space/time diagrams are also often employed as the user display in concurrent program debugging tools.

The Los Alamos parallel debugging system uses a text based *time-process* display, and Idd uses a graphic display. Both of these, however, rely on an accurate global real-time clock (impractical in most systems).

Fig. 33 shows a partial order event tracer (POET) display 3300. The partial order event tracer (POET) system supports several different languages and run-time environments, including Hermes, a high-level interpreted language for distributed systems, and Java. With reference to Fig. 33, POET display 3300 distinguishes among several types of events by shapes, shading, and alignment of corresponding message lines.

A Distributed Program Debugger (DPD) is based on an Remote Execution Manager (REM) framework. The REM framework is a set of servers on interconnected Unix machines, where each server is a Unix user-level process. Processes executing in this framework can create and communicate with processes elsewhere in the network as if they were all on the same machine. DPD uses

space/time displays for debugging communication only, and it relies on separate source-level debuggers for individual processes.

2. Abstraction in event-based debugging

5 Simple space/time displays can be used to present programmers with a wealth of information about distributed executions. Typically however, space/time diagrams are too abstract to be an ultimate debugging solution. Space/ time diagrams show high-level events and message traffic, but they do not support designer interaction with the source code. On the other hand, simple space/time diagrams may sometimes
10 have too much detail. Space/time diagrams display each distinct low-level message that contributes to a high-level transaction without support for abstracting the transaction.

Fig. 34 is a space/time diagram 3400 having a first compound event 3402 and a second compound event 3404. With reference to Fig. 34, even though a pair of
15 primitive events are either causally related or concurrent, first and second compound events 3402 and 3404, or any other pair of compound events, might be neither causally related nor concurrent. Abstraction is typically applied across two dimensions—events and processes—to aid in the task of debugging distributed software. Event abstraction represents sequences of events as single entities. A group
20 of events may occasionally have a specific semantic meaning that is difficult to recognize, much as streams of characters can have a meaning that is difficult to interpret without proper spacing and punctuation. Event abstraction can in some circumstances complicate the relationships between events.

Event abstraction can be applied in one of three ways: filtering, clustering, and
25 interpretation. With event filtering, a programmer describes event types that the debugger should ignore, which are then hidden from view. With clustering, the debugger collects a number of events and presents the group as a single event. With interpretation, the debugger parses the event stream for event sequences with specific semantic meaning and presents them to a programmer.

Process abstraction is usually applied only as hierarchical clustering. The remainder of this section discusses these specific event and process abstraction approaches.

A. Event filtering and clustering

5

Event filtering and clustering are techniques used to hide events from a designer and thereby reduce clutter. Event filters exclude selected events from being tracked in event-based debugging techniques. In most cases, this filtering is implicit and can not be modified without changing the source code because the source code
10 being debugged is designed to report only certain events to the debugger. When deployed, the code will report all such events to the tool. This approach is employed in both DPD and POET, although some events may be filtered from the display at a later time.

An event cluster is a group of events represented as a single event. The
15 placement of an event in a cluster is based on simple parameters, such as virtual time bounds and process groups. Event clusters can have causal ambiguities. For example, one cluster may contain events that causally precede events in a second cluster, while other events causally follow certain events in the second cluster.

Fig. 35 shows a POET display 3500 involving a first convex event cluster
20 3502 and a second convex event cluster 3504. POET uses a virtual-time-based clustering technique that represents convex event clusters as single abstract events. A convex event cluster is a set of event instances, C , such that for events

$$a, b, c \in E \text{ with } a, c \in C, a \sim b \wedge b \sim c \Rightarrow b \in C$$

25

Convex event clusters, unlike generic event clusters, cannot overlap.

B. Event interpretation

The third technique for applying event abstraction is interpretation, also referred to as behavioral abstraction. Both terms describe techniques that use debugging tools to interpret the behavior represented by sequences of events and present results to a designer. Most approaches to behavioral abstraction let a designer describe sequences of events using expressions, and the tools recognize the sequence of events through a combination of customized finite automata followed by explicit checks. Typically, matched expressions generate new events.

1. Event description language (EDL)

One of the earliest behavioral abstraction technique was Bates's event description language (EDL), where event streams are pattern-matched using shuffle automata. A match produces a new event that can, in turn, be part of another pattern. Essentially, abstract events are hierarchical and are built bottom up.

This approach can recognize event patterns that contain concurrent events. There are however several weaknesses in this approach. First, shuffle automata match events from a linear stream, which is subject to a strong observational bias. In addition, even if the stream constitutes a valid observation, interleaving may cause false intermediates between an event and its immediate successor. Finally, concurrent events appear to occur in some specific order.

Bates partially compensates for these problems in three ways. First, all intermediates between two recognized events are ignored—hence, false intermediates are skipped. Unfortunately, true intermediates are also skipped, making error detection difficult. Second, the shuffle operator, Δ , is used to identify matches with concurrent events. Unfortunately, shuffle recognizes events that occur in any order, regardless of whether they are truly ordered in the corresponding execution. For example, $e_1 \Delta e_2$ can match with either $e_1 < e_2$ or $e_2 < e_1$ in the event stream, but this means the actual matches could be: $e_1 \sim e_2$, $e_2 \sim e_1$, in addition to the $e_1 \parallel e_2$ that the programmer intended to match. Third, the programmer can prescribe explicit checks

to be performed on each match before asserting the results. However, the checks allowed do not include causality or concurrence checks.

2. Chain expressions

Chain expressions, used in the Ariadne parallel debugger, are an alternate way to describe distributed behavior patterns that have both causality and concurrence. These behavioral descriptions are based on chains of events (abstract sequences not bound to processes), p-chains (chains bound to processes), and pt-chains (composed p-chains). The syntax for describing chain expressions is fairly simple, with $\langle a \ b \rangle$ representing two causally related events and $|[a \ b]|$ representing two concurrent events.

The recognition algorithm has two functions. First, recognizing the appropriate event sequence from a linear stream (using an NFA). Second, checking the relationships between specific events

For example, when looking for sequences that match the expression $\langle |[a \ b]| \ c \rangle$ (viz., a and b are concurrent, and both causally precede c), Ariadne will find the sequence a b c and then verify the relationships among them. Unfortunately, the fact that sequences are picked in order from a linear stream before relationships are checked can cause certain matches to be missed. For example, $|[a \ b]|$ and $|[b \ a]|$ should have the same meaning, but they do not cause identical matches. This is because Ariadne uses NFAs as the first stage in event abstraction. In the totally ordered stream to which an NFA responds, either a will precede b, preventing the NFA for the second expression from recognizing the string, or b will precede a, preventing the NFA for the first expression from recognizing the string.

3. Distributed abstraction

The behavioral abstraction techniques described so far rely on centralized abstraction facilities. These facilities can be distributed, as well. The BEE (Basis for distributed Event Environments) project is a distributed, hierarchical, event-collection system, with debugging clients located with each process.

Fig. 36 show a Basis for distributed Event Environments (BEE) abstraction facility 3600 for a single client. With reference to Fig. 36, event interpretation is performed at several levels. The first is an event sensor 3602, inserted into the source of the program under test and invoked whenever a primitive event occurs during execution. The next level is an event generator 3604, where information—including timestamps and process identifiers—is attached to each event. Event generator 3604 uses an event table 3606 to determine whether events should be passed to an event handler 3608 or simply dropped. Event handler 3608 manages event table 3606 within event generator 3604. Event handler 3608 filters and collects events and routes them to appropriate event interpreters (not shown). Event interpreters (not shown) gather events from a number of clients (not shown) and aggregate them for presentation to a programmer. Clients and their related event interpreters are placed together in groups managed by an event manager (not shown). A weakness of this technique is that it does not specifically track causality. Instead, this technique relies on the real-time stamps attached to specific primitive or abstract events. However, as discussed above these timestamps are not able to characterize causality.

C. Process clustering

Most distributed computing environments feature flat process structures, with few formally stated relationships among processes. Automatic process clustering tools can partially reverse-engineer a hierarchical structure to help remove spurious information from a debugger's view. Intuitively, a good cluster hierarchy should reveal, at the top level, high-level system behavior, and the resolution should improve proportionally with the number of processes exposed. A poor cluster hierarchy would show very little at the top level and would require a programmer to descend several hierarchical levels before getting even a rough idea about system behavior. Process clustering tools attempt to identify common interaction patterns—such as client-server, master-slave, complex server, layered system, and so forth. When these patterns are identified, the participants are clustered together. Clusters can then serve as participants in interaction patterns to be further clustered. These cluster hierarchies are strictly trees, as shown in Fig. 37, which depicts a hierarchical construction of

process clusters 3700. With reference to Fig. 37, a square node 3702 represents a process (not shown) and a round node 3704 represents a process cluster (not shown).

5 Programmers can choose a debugging focus, in which they specify the aspects and detail levels they want to use to observe an execution. With reference to Fig. 37, a representative debugging focus that includes nodes I, J, E, F, G and H is shown. One drawback of this approach is that when a parent cluster is in focus, none of its children can be. For example, if we wanted to look at process K in detail, we would also need to expose at least as much detail for processes E and L and process cluster D.

10 Each process usually participates in many types of interactions with other processes. Therefore, the abstraction tools must heuristically decide between several options. These decisions have a substantial impact on the quality of a cluster hierarchy. In "Abstract Behaviour of Distributed Executions with Applications to Visualization," PhD. thesis, Technische Hochschule Darmstadt, Darmstadt, Germany, 15 May 1994, by T. Kunz, the author evaluates the quality of his tool by measuring the cohesion, which though expressed quantitatively is actually a qualitative measurement (the higher the better) within a cluster and the coupling, a qualitative measure of the information clusters must know about each other (the higher the worse), between clusters. For a cluster P of m processes, cohesion is quantified by:

20

$$\text{Cohesion}(P) = \frac{\sum_{i < j} \text{Sim}_f(p_i, p_j)}{m(m-1)/2}$$

where $\text{Sim}_f(P_1, P_2)$ is a similarity metric that equals:

25

$$\text{Sim}_f = \frac{A\langle \hat{C}_{P_1} | \hat{C}_{P_2} \rangle}{\| \hat{C}_{P_1} \| \cdot \| \hat{C}_{P_2} \|}$$

Here, $\langle \hat{a} | \hat{b} \rangle$ denotes the scalar product of vectors \hat{a} and \hat{b} , and $\|\hat{a}\|$ denotes the magnitude of vector \hat{a} . C_{P_1} and C_{P_2} are process characteristic vectors—in them, each element contains a value between 0 and 1 that indicates how strongly a particular characteristic manifests itself in each process. Characteristics can include keywords, type names, function references, etc. A is a value that equals 1 if any of the following apply:

- P_1 and P_1 are instantiations of the same source.
- P_1 and P_2 are unique instantiations of their own source.
- P_1 and P_2 communicate with each other.

A equals 0 if none of these is true (e.g., P_1 and P_2 are non-unique instantiations of separate source that do not communicate with each other). Coupling is quantified by:

$$\text{Coupling}(P) = \frac{\sum_{ij} \text{Sim}_f(p_i, q_j)}{mn}$$

15

where $q_j \in Q$, Q is the complement of P , and $n = |Q|$. The quality of a cluster is quantified as its Coupling minus its Cohesion. In many cases, these metrics match many of the characteristics that intuitively differentiate good and poor clusters, as shown in Fig. 38A, B, and C. With reference to Fig. 38A and C, Cohesion is high where clusters correspond to heavy communication and where clusters correspond to processes instantiated from the same source code. Coupling is shown to be low in each of the above cases. With reference to Fig. 38B Coupling is high when clusters do not correspond to heavily communicating processes or to instances of the same source code. It is not clear, however, that the cluster in Fig. 38C should be assigned the same quality value as the cluster in Fig. 38A. Using these metrics, Kunz achieved qualities of between :15 and :31 for his clustering techniques. However, it is hard to tell what this means in terms of cluster usefulness.

3. State-based debugging

State-based debugging techniques focus on the state of the system and the state changes caused by events, rather than on events themselves. The familiar source-level debugger for sequential program debugging is state-based. This source-level debugger lets designers set breakpoints in the execution of a program, enabling them to investigate the state left by the execution to that point. This source-level debugger also lets programmers step through a program's execution and view changes in state caused by each step.

Concurrent systems have no unique meaning for an instant in execution time. Stopping or single stepping the whole system can unintentionally, but substantially, change the nature of interactions between processes.

A. Consistent cuts and global state

In distributed event-based debugging, the concept of causality is typically of such importance that little of value can be discussed without a firm understanding of causality and its implications. In distributed state-based debugging, the concept of a global instant in time is equally important.

Here again, it may seem intuitive to consider real-time instants as the global instants of interest. However, just as determining the real-time order of events is not practical or even particularly useful, finding accurate real-time instants makes little sense. Instead, a global instant is represented by a consistent cut. A consistent cut is a cut of an event dependency graph representing an execution that: a) intersects each process exactly once; and b) points all dependencies crossing the cut in the same direction. Like real-time instants, consistent cuts have both a past and a future. These are the subgraphs on each side of the cut.

Fig. 39, shows that consistent cuts can be represented as a jagged line across the space/time diagram that meets the above requirements. With reference to Fig. 39, a space/ time graph 3900 is shown having a first cut 3902 and a second cut 3904. All events to the left of either first cut 3902 or second cut 3904 are in the past of each cut, and all events to the right are in the future of each cut, respectively. First cut

3902 is a consistent cut because no message travel from the future to the past. Second cut 3904 however is not consistent because a message 3906 travels from the future to the past.

Fig. 40A, B, and C show that a distributed execution shown in a space/time diagram 4000 can be represented by a lattice of consistent cuts 4002, where \top is the start of the execution, and \perp is system termination. With reference to Fig. 40A, B, and C, lattice of consistent cuts 4002 represents the global statespace traversed by a single execution. Since lattice of consistent cuts 4002's size is on the order of $|E|^{|\mathcal{P}|}$ it, unlike space/time diagrams, are never actually constructed. In the remainder of this chapter, to describe properties of consistent cut lattices, the symbol \xrightarrow{l} relates cuts such that one immediately precedes the other, and \sim relates cuts between which there is a path.

B. Single stepping in a distributed environment

Controlled stepping, or single stepping, through regions of an execution can help with an analysis of system behavior. The programmer can examine changes in state at the completion of each step to get a better understanding of system control flow. Coherent single stepping for a distributed system requires steps to align with a path through a normal execution's consistent cut lattice.

DPD works with standard single-process debuggers (called client debuggers), such as DBX, GDB, etc. Programmers can use these tools to set source-level breakpoints and single step through individual process executions. However, doing so leaves the other processes executing during each step, which can yield unrealistic executions.

Zernic gives a simple procedure for single stepping using a post-mortem traversal of a consistent cut lattice. At each point in the step process, there are two disjoint sets of events: the past set, or events that have already been encountered by the stepping tool, and the future set, or those that have yet to be encountered. To perform a step, the debugger chooses an event, e_i , from the future such that any

events it depends on are already in the past, i.e., there are no future events, e_f , such that $e_f \rightsquigarrow e_i$. This ensures that the step proceeds between two consistent cuts related by \xrightarrow{l} . The debugger moves this single event to the past, performing any necessary actions.

- 5 To allow more types of steps, POET's support for single stepping uses three disjoint sets: executed, ready, and non-ready. The executed set is identical to the past set in "Using Visualization Tools to Understand Concurrency," by D. Zernik, M. Snir, and D. Malki, *IEEE Software* 9, 3 (1992), pp. 87-92. The ready set contains all events that are fully enabled by events in the future, and the contents of the non-ready
- 10 set have some enabling events in either the ready or non-ready sets. Using these sets, it is possible to perform three different types of step: global-step, step-over, and step-in. Global-step and step-over may progress between two consistent cuts not related \xrightarrow{l} (i.e., there may be several intermediate cuts between the step cuts).

- A global-step is performed by moving all events from the ready set into the
- 15 past. Afterwards, the debugger must move to the ready set all events in the non-ready set whose dependencies are in the executed set. A global-step is useful when the programmer wants information about a system execution without having to look at any process in detail.

- The step-over procedure considers a local, or single-process, projection of the
- 20 ready and non-ready sets. To perform a step, it moves the earliest event from the local projections into the executed set and executes through events on the other processes until the next event in the projection is ready. This ensures that the process in focus will always have an event ready to execute in the step that follows.

- Step-in is another type of local step. Unlike step-over, step-in does not
- 25 advance the system at the completion of the step; instead, the system advance is considered to be a second step. Fig. 41A, B, C, and D show a space/time diagram before a step 4100 and a resulting space/time diagram after performing a global-step 4102, a step-over 4104, and a step-in 4106.

C. Runtime consistent cut algorithms

It is occasionally necessary to capture consistent cuts at runtime. To do so, each process performs some type of cut action (e.g., state saving). This can be done
5 with barrier synchronization, which erects a temporal barrier that no process can pass until all processes arrive. Any cut taken immediately before, or immediately after, the barrier is consistent. However, with barrier synchronization, some processes may have a long wait before the final process arrives.

A more proactive technique is to use a process called the cutinitiator to send
10 perform-cut messages to all other system processes. Upon receiving a perform-cut message, a process performs its cut action, sends a cut-finished message to the initiator, and then suspends itself. After the cut initiator receives cut-finished messages from all processes, it sends each of them a message to resume computation.

The cut obtained by this algorithm is consistent: no process is allowed to send
15 any messages from the time it performs its own cut action until all processes have completed the cut. This means that no post-cut messages can be received by processes that have yet to perform their own cut action. This algorithm has the undesirable characteristic of stopping the system for the duration of the cut. The following algorithms differ in that they allow some processing to continue.

20 1. Chandy-Lamport algorithm

The Chandy-Lamport algorithm does not require the system to be stopped. Once again, the cut starts when a cut initiator sends perform-cut messages to all of the processes. When a process receives a perform-cut message, it stops all work, performs its cut action, and then sends a mark on each of its outgoing channels; a
25 mark is a special message that tells its recipient to perform a cut action before reading the next message from the channel. When all marks have been sent, the process is free to continue computation. If the recipient has already performed the cut action when it receives a mark, it can continue working as normal.

Each cut request and each mark associated with a particular cut are labeled
30 with a cut identifier, such as the process ID of the cut initiator and an integer. This

lets a process distinguish between marks for cuts it has already performed and marks for cuts it has yet to perform.

2. Color-based algorithms

The Chandy-Lamport algorithm works only for FIFO (First In First Out) channels. If a channel is non-FIFO (*e.g.*, UDP), a post-cut message may outrun the mark and be inconsistently received before the recipient is even aware of the cut, *i.e.*, it is received in the cut's past. The remedy to this situation is a color-based algorithm. Two such algorithms are discussed below.

The first is called the two-color, or red-white, algorithm. With this algorithm, information about the cut state is transferred with each message. Each process in the system has a color. Processes not currently involved in a consistent cut are white, and all messages transmitted are given a white tag. Again, there is a cut initiator that sends perform-cut messages to all system processes. When a process receives this request, it halts, performs the cut action, and changes its color to red. From this point on, all messages transmitted are tagged with red to inform the recipients that a cut has occurred.

Any process can accept a white message without consequence, but when a white process receives a red message, it must perform its cut action before accepting the message. Essentially, white processes treat red messages as cut requests. Red processes can accept red messages at any time, without consequence.

A disadvantage of the two-color algorithm is that the system must reset all of the processes back to white after they have completed their cut action. After switching back, each process must treat red messages as if they were white until they are all flushed from the previous cut. After this, each process knows that the next red message it receives signals the next consistent cut.

This problem is addressed by the three-color algorithm, which resembles the two-color algorithm in that every process changes color after performing a cut; it differs in that every change in color represents a cut. For colors zero through two, if a process with the color c receives a message with the color $(c - 1) \bmod 3$, it

registers this as a message-in-flight (see below). On the other hand, if it receives a message with the color $(c + 1) \bmod 3$, it must perform its cut action and switch color to $(c + 1) \bmod 3$ before receiving the message. Of course, this can now be generalized to n-color algorithms, but three colors are usually sufficient.

5 Programmers may need to know about messages transmitted across the cut, or messages-in-flight. In the two-color algorithm, messages-in-flight are simply white messages received by red processes. These can all be recorded locally, or the recipient can report them to the cut initiator. In the latter case, each red process simply sends the initiator a record of any white messages received.

10 It is not safe to switch from red to white in the two-color algorithm until the last message-in-flight has been received. This can be detected by associating a counter with each process. A process increments its counter for each message sent and decrements it for each message received. When the value of this counter is sent to the initiator at the start of each process's cut action, the initiator can use the total
15 value to determine the total number of messages-in-flight. The initiator simply decrements this count for each message-in-flight notification it receives.

D. State recovery—rollback and replay

 Since distributed executions tend to be non-deterministic, it is often difficult to
20 reproduce bugs that occur during individual executions. To do so, most distributed debuggers contain a rollback facility that returns the system to a previous state. For this to be feasible, all processes in the system must occasionally save their state. This is called checkpointing the system. Checkpoints do not have to save the entire state of the system. It is sufficient to save only the changes since the last checkpoint.
25 However, such incremental checkpointing can prolong recovery.

 DPD makes use of the UNIX fork system call to perform checkpointing for later rollback. When fork is called, it makes an exact copy of the calling process, including all current state. In the DPD checkpoint facility, the newly forked process is suspended and indexed. Rollback suspends the active process and resumes an
30 indexed process. The problem with this approach is that it can quickly consume all

system memory, especially if checkpointing occurs too frequently. DPD's solution is to let the programmer choose the checkpoint frequency through use of a slider in its GUI.

Processes must sometimes be returned to states that were not specifically saved. In this case, the debugger must do additional work to advance the system to the desired point. This is called replay and is performed using event trace information to guide an execution of the system. In replay, the debugger chooses an enabled process (*i.e.*, one whose next event has no pending causal requirements) and executes it, using the event trace to determine where the process needs to block for a message that may have arrived asynchronously in the original execution. When the process blocks, the debugger chooses the next enabled process and continues from there. In this way, a replay is causally identical to the original execution.

Checkpoints must be used in a way that prevents domino effects. The domino effect occurs when rollbacks force processes to restore more than one state. Domino effects can roll the system back to the starting point. Fig. 42 shows a space time diagram 4200 for a system that is subject to the domino effect during rollback. With reference to Fig. 42, if the system requests a rollback to checkpoint c_3 4202 of process P_3 4204, all processes in the system must roll back to c_1 (*i.e.*, roll back to $P_3.c_2$ 4206 requires a roll back to $P_2.c_2$ 4208, which requires a roll back to $P_1.c_2$ 4210, which requires a roll back to $P_3.c_1$ 4212, which requires a roll back to $P_2.c_1$ 4214, which requires a final roll back to $P_1.c_1$ 4216). The problem is caused by causal overlaps between message transfers and checkpoints. Performing checkpoints only at consistent cuts avoids a domino effect.

E. Global state predicates

The ability to detect the truth value of predicates on global state yields much leverage when debugging distributed systems. This technique lets programmers raise flags when global assertions fail, set global breakpoints, and monitor interesting aspects of an execution. Global predicates are those whose truth value depends on the state maintained by several processes. They are typically denoted with the symbol Φ .

Some examples include $(\sum_i c_i > 20)$ and $(c_1 < 20 \wedge c_2 > 5)$, where c_i is some variable in process P_i that stores positive integers. In the worst case (such as when $(\sum_i c_i > 20)$ is false for an entire execution), it may be necessary to get the value of all such variables in all consistent cuts. In the following discussion, we use the notation C_a

5 $\models \Phi$ to indicate that Φ is true in consistent cut C_a .

At this point, it is useful to introduce branching time temporal logic. Branching time temporal logic is predicate logic with temporal quantifiers, P, F, G, H, A, and E. $P\Phi$ is true in the present if Φ was true at some point in the past; $F\Phi$ is true in the present if Φ will be true at some point in the future; $G\Phi$ is true in the

10 present if Φ will be true at every moment in the future; and $H\Phi$ is true in the present if Φ was true at every moment of the past. Notice that $G\Phi$ is the same as $\neg F\neg\Phi$, and $H\Phi$ is the same as $\neg P\neg\Phi$.

Since global time passage in distributed systems is marked by a partially ordered consistent cut lattice rather than by a totally ordered stream, we need the

15 quantifiers A, which precedes a predicate that is true on all paths, and E, which precedes a predicate that is true on at least one path. So, $AF\Phi$ is true in the consistent cut representing the present if Φ is true at least once on all paths in the lattice leaving this cut. $EP\Phi$ is true in the consistent cut representing the present if Φ is true on at least one path leading to this cut.

20 A monotonic global predicate is a predicate Φ such that $C_a \models \Phi \Rightarrow C_a \models AG\Phi$. A monotonic global predicate is one that remains true after becoming true. An unstable global predicate, on the other hand, is a predicate Φ such that $C_a \models \Phi \Rightarrow C_a \models EG\neg\Phi$. An unstable global predicate is one that may become false after becoming true.

25 1. Detecting monotonic global predicates

Monotonic predicates can be detected any time after becoming true. One algorithm is to occasionally take consistent cuts and evaluate the predicate at each. In fact, it is not necessary to use consistent cuts, since any transverse cut whose future is

a subset of the future of the consistent cut in which the predicate first became true will also show the predicate true.

2. Detecting unstable global predicates

Detecting arbitrary unstable global predicates can take at worst $|E|^{|\mathcal{P}|}$ time, where $|E|^{|\mathcal{P}|}$ is the size of an execution's consistent cut lattice, where $|E|$ is the number of events in the execution, and $|\mathcal{P}|$ is the number of processes. This is so, because it may be necessary to test for the predicate in every possible consistent cut. However, there are a few special circumstances that allow $|E|$ time algorithms.

Some unstable global predicates are true on only a few paths through the consistent cut lattice, while others are true on all paths. Cooper and Marzullo describe predicate qualifiers definitely Φ for predicates that are true on all paths (*i.e.*, $\top \models A \models \Phi$) and possibly Φ for those that are true on at least one path (*i.e.*, $\top \models \exists E \models \Phi$).

The detection of possibly Φ for weak conjunctive predicates, or global predicates that can be expressed as conjunctions of local predicates, is $\phi(|E|)$. The algorithm for this is to walk a path through the consistent cut lattice that aligns with a single process, P_i , until either: (1) the process's component of Φ is true, or (2) there is no way to proceed without diverging from P_i . In either case, the target process is switched, and the walk continued. This algorithm continues until it reaches a state where all components of the predicate are true, or until it reaches \perp . In this way, if there are any consistent cuts where all parts of the predicate simultaneously hold, the algorithm will encounter at least one.

Detection of possibly Φ for weak disjunctive predicates, or global predicates that can be expressed as disjunctions of local predicates, is also $\phi(|E|)$; it is the same algorithm as above, except it halts at the first node where *any* component is true. However, weak conjunctive and disjunctive predicates constitute only a small portion of the types of predicates that could be useful in debugging distributed systems.

3. Conclusions

Complicating the debugging of heterogenous embedded systems are designs comprised of concurrent and distributed processes. Most of the difficulty in debugging distributed systems results from concurrent processes with globally
5 unscheduled and frequently asynchronous interactions. Multiple executions of a system can produce wildly varying results—even if they are based on identical inputs. The two main debugging approaches for these systems are event based and state based.

10 Event-based approaches are monitoring approaches. Events are presented to a designer in partially ordered event displays, called space/time displays. These are particularly good at showing inter-process communication over time. They can provide a designer with large amounts of information in a relatively small amount of space.

15 State-based approaches focus locally on the state of individual processes or globally on the state of the system. Designers can observe individual system states, set watches for specific global predicates, step through executions, and set breakpoints based on global state predicates. These approaches deal largely with snapshots, considering temporal aspects only as differences between snapshots.

20 As distributed systems increase in size and complexity, the sheer volume of events generated during an execution grows to a point where it is exceedingly difficult for designers to correctly identify aspects of the execution that may be relevant in locating a bug. For distributed system debugging techniques to scale to larger and faster systems, behavioral abstraction will typically become a necessity to help designers identify and interpret complicated behavioral sequences in a system
25 execution. Finally, embedded systems must execute in a separate environment from the one in which they were designed and embedded systems may also run for long periods of time without clear stopping points. Debugging them requires probes to report debugging information to a designer during the execution. These probes inevitably alter system behavior, which can mask existing bugs or create new bugs

that are not present in the uninstrumented system. While it is not possible to completely avoid these probe effects, they can be minimized through careful placement, or masked through permanent placement.

Debugging by Cooperative Execution in Coordination-Centric Designs

- 5 Fig. 43 is an example of a coordination-centric approach to the debugging of a distributed software environment 4300, in accordance with the present invention. With reference to Fig. 43, distributed software environment 4300, containing several processing elements 4302 connected so that information can be exchanged, is connected through a communication channel 4304 to a debugging host 4306.
- 10 Distributed software environment 4300 can be connected either directly or indirectly to debugging host 4306.

- Debugging, in a coordination-centric approach, is performed using "cooperative execution." Cooperative execution refers to simultaneously executing a distributed software environment 4300 and simulating distributed software
- 15 environment 4300 on debugging host 4306 based on event traces from distributed software environment 4300. Debugging host 4306 may be a general-purpose workstation. Distributed software environments are likely to have several processing elements 4302, but only a few of these have resources that let them connect directly to debugging host 4306. Those that do not have these resources, either have an indirect
- 20 connection to the debugging host or are opaque to debugging.

1. Direct Connection

- Events can be recorded locally on distributed software environment 4300, and event records can be transferred directly to debugging host 4306. Fig. 44 is a detailed view of a direct connection between a primary processing element 4400 of distributed
- 25 software environment 4300 and debugging host 4306. In a direct connection, primary processing element 4400 contains a software program 4401, which has at least two software components 4402, a runtime system 4404, a coordinator 4406, an interface (not shown) having a port (not shown), and a primary runtime debugging architecture

4408. Software program 4401 generates an event record in response to a selected event.

Runtime system 4404 collects the events from software components 4402 and transfers them to primary runtime debugging architecture 4408. Primary runtime debugging architecture 4408 contains a time stamper 4410, a causality stamper 4412, a primary uplink component 4414, and a primary transfer component 4416 coupled to primary uplink component 4414. Time stamper 4410 provides a time stamp to the event record generated by software program 4401, while causality stamper 4412 provides an identification of a cause of the event associated with the corresponding event record. Primary uplink component 4414 of primary runtime debugging architecture 4408 facilitates communication through communication channel 4304 between primary processing element 4400 and debugging host 4306, while primary transfer component 4416 collects and facilitates the transfer of the time-stamped and causality-stamped event record from primary processing element 4400 to an event queue 4418 on debugging host 4306.

Debugging host 4306 operates on the events—simulating the activity of distributed software environment 4300 and letting the designer navigate the execution. The bulk of debugging support is thus on debugging host 4306, which reduces the probe effect in distributed software environment 4300. All debugging is based on a combination of system state and event records, which must be collected and properly annotated to record causality.

With a direct connection, primary processing element's runtime system 4404 asynchronously dumps all events onto debugging host 4306. Debugging host 4306 queues all events and writes them to disk. All debugging activity is based on post-processing on these events.

2. Indirect Connection

Fig. 45 is a detailed view of an indirect connection from primary processing element 4400 to debugging host 4306. In an indirect connection, primary runtime debugging architecture 4408 contains time stamper 4410, causality stamper 4412, and

primary uplink component 4414, whereas an intermediate processing element 4500 contains an intermediate runtime debugging architecture 4502, having an intermediate transfer component 4504 coupled to an intermediate uplink component 4506.

For primary processing element 4400 with an indirect connection to debugging host 4306, event records are routed to intermediate processing element 4500, which collects and forwards them, along with some event records of its own, to debugging host 4306. Primary uplink component 4414 facilitates communication, along communication channel 4304, between primary processing element 4400 and intermediate processing element 4500. Intermediate uplink component 4506 facilitates communication, along communication channel 4304, between intermediate processing element 4500 and debugging host 4306, while intermediate transfer component 4504 collects and facilitates the transfer of the event record from intermediate processing element 4500 to event queue 4418 in debugging host 4306.

3. Storage-Based Event Collection

If it is not possible or practical to connect distributed software environment 4300 to debugging host 4306, it may still be possible to employ cooperative execution in a post-mortem fashion. Fig. 46 depicts capturing event records in a flash memory for post-mortem distributed debugging. In the present invention, runtime system 4404 collects and transfers the events to primary runtime debugging architecture 4408, which includes time stamper 4410, causality stamper 4412, and a flash driver 4600. After the events have been time-stamped and causality-stamped, flash driver 4600 sends the events to a flash memory 4602 for storage and post-mortem debugging.

As long as a global name space is used, each processing element 4302 can store events separately, and the complete global event trace can be reconstructed later. Unlike connection-based collection, there is a definite limit on the amount of storage space, and it is essential to use it wisely. Fig. 47 shows how flash memory 4602 can be allocated to ensure that the entire context of a system crash can be reconstructed.

4. Event Recording Call

The runtime system of a processing element records events from the execution, applies timestamps and causality pointers, and forwards the records to a communication channel. Any event made visible to the runtime system, either by
5 being identified as an explicit event in the model or by *eventRecord* calls at runtime, is a candidate for recording. The *eventRecord* function can be implemented either to transfer event records to a runtime simulation host or to record them into storage.

To build concurrent event traces, distributed software environment 4300 must be instrumented to produce appropriate event records, and media must be applied to
10 store them. Unlike traditional approaches, instrumentation must include information about how different software components 4402 interact. Instrumented distributed code can change the execution characteristics of distributed software environment 4300 being debugged (recall the *probe effect*).

To collect event traces, the code itself must be well instrumented. This
15 instrumentation can be automatically applied to the code by the runtime synthesis engine. A simple form of instrumentation involves inserting event recording calls at each significant source line (such as *message sent* or *variable changed*) and can include the values sent or the values to which the variable was changed. Over-parameterization of event records can create significant delays in simulated systems
20 and can cause a probe effect in physical systems, since the time for packaging parameters can be significant. For the most part, static instrumentation code is required for activities that do not appear as events to runtime system 4404, since runtime system 4404 can deal with explicit events. The *eventRecord* method assigns a unique identifier to the event record, assigns all appropriate timestamps, and identifies
25 implicit causality. See Listing 2 below for an example of an event recording call:

Listing 2 - Instrumented Event Record

```
cl.p2.send (x + 3);  
eventRecord ("cl.p2.send"); // synthesized
```


As with other aspects of the system, a designer does not know in advance exactly how much detail will be needed, and detail requirements may change over time. This means that the principles of selective focus should be applied to event
5 collection, as well. While static instrumentation cannot be altered during system execution, it can be designed to be disabled by the debugging system during execution.

Runtime debugging support in the physical distributed software environment should be as lightweight as possible to avoid incurring any further probe effect. One
10 of the goals in debugging a physical system is to ensure that as little of the debugging support as possible is placed on the systems themselves.

5. Precomputed Event Sequences

Events often occur in specific, predetermined, partially ordered sequences. Therefore, instead of sending every single event in the sequence, it is possible to send
15 a token representing the sequence and let debugging host 4306 expand it. This is an application of selective focus to event collection. Since the lowest detail levels are completely determined, there is nothing lost when operating at these higher levels of abstraction.

For example, with the RPC protocol, it may be necessary to indicate only the
20 start and finish of a transaction. If a bug causes the transaction to fail, simply knowing the identity of that particular transaction can be a useful starting point for further evaluation in simulation.

6. Execution Shadowing

Since individual traces only include events, the simulator must keep track of
25 all other aspects of an execution. The host side of the debugger tracks system state through state change events. This lets the debugger track decisions made by the distributed system and is called execution shadowing.

7. Interpolation of Events from Bus Traces

Events that are visible on a system bus can be recorded by observing bus signals. This requires hardware support for bus trace collection. In this case, events of a far more primitive nature must be dealt with (*e.g.*, transitions between states for physical wires and very low-level actions).

Fig. 48 shows an example of low-level behavioral recognition for a simple I²C protocol. At the lowest level are signal transitions. These are translated into sub-low-level events, which are then fed into a language recognizer to generate low-level debugging events.

The main problems with this approach involve bandwidth. Potentially large numbers of events must be analyzed to present a reasonable view of what happens at the lower levels. This can either (1) require the hardware to save a great deal of information and then parse it later for content or (2) cause a great deal of computation overhead at runtime, which slows down the system.

It will be obvious to those having skill in the art that many changes may be made to the details of the above-described embodiments of this invention without departing from the underlying principles thereof. The scope of the present invention should, therefore, be determined only by the following claims.